

# HIGH-PERFORMANCE PARALLEL INTERFACE - Scheduled Transfer (HIPPI-ST)

May 8, 1997

Secretariat:

Information Technology Industry Council (ITI)

**ABSTRACT:** This standard specifies a data transfer protocol that uses small control messages to pre-arrange data movement. Buffers are allocated at each end before the data transmission, allowing full-rate, non-congesting data flow between the end devices. The control and data may use different physical media, or may share a single physical medium. Procedures are provided for moving data over HIPPI and other media.

**NOTE:**

*This is an internal working document of X3T11, a Technical Committee of Accredited Standards Committee X3. As such, this is not a completed standard. The contents are actively being modified by X3T11. This document is made available for review and comment only. For current information on the status of this document contact the individuals shown below:*

**POINTS OF CONTACT:**

Roger Cummings (X3T11 Chairman)  
Distributed Processing Technology  
140 Candace Drive  
Maitland, FL 32751  
(407) 830-5522 x348, Fax: (407) 260-5366  
E-mail: cummings\_roger@dpt.com

Ed Grivna (X3T11 Vice-Chairman)  
Cypress Semiconductor  
2401 East 86th Street  
Bloomington, MN 55425  
(612) 851-5200, Fax: (612) 851-5087  
E-mail: elg@cypress.com

Don Tolmie (HIPPI-ST Technical Editor)  
Los Alamos National Laboratory  
CIC-5, MS-B255  
Los Alamos, NM 87545  
(505) 667-5502, Fax: (505) 665-7793  
E-mail: det@lanl.gov

## Comments on Rev 0.6

This is a preliminary document undergoing lots of changes. Many of the additions are just place holders, or are put there to stimulate discussion. Hence, do not assume that the items herein are correct, or final – everything is subject to change. This page tries to outline where we are; what has been discussed and semi-approved, and what has been added or changed recently and deserves your special attention. This summary relates to changes since the previous revision. Also, previous open issues are outlined with a single box, new open issues ones are marked with a double bar on the left edge of the box.

Changes are marked with margin bars so that changed paragraphs are easily found, and then highlights mark the specific changes. The list below just describes the major changes, for detail changes please compare this revision to the previous revision.

Please help us in this development process by sending comments, corrections, and suggestions to the Technical Editor, Don Tolmie, of the Los Alamos National Laboratory, at det@lanl.gov. If you would like to address the whole group working on this document, send the comment(s) to hippy@network.com.

1. In 2, deleted the normative reference for HIPPI-SC.
2. Did a global change of "set up" back to "setup", and then put several back to "set up". The difference is that "set up" is used as a verb, and "setup" as a noun. These changes are not marked in the document.
3. In 3.1, added a definition for "Opaque data".
4. In 3.3, deleted the acronym for MPI since it wasn't used in the document. Changed "SR State\_Response" to "RSR Request\_State\_Response". Added "RTRR Request\_To\_Receive\_Response".
5. Throughout the document, changed "State\_Response" to "Request\_State\_Response". These changes are highlighted, but do not have margin bars.
6. In 4.3.2, added the last sentence of the 3rd paragraph about Port = x'0000' being valid only for Request\_Port Operations.
7. In 4.3.4, changed the maximum size of a Buffer from  $2^{32}$  to  $2^{64}$  (a very large number).
8. In 4.3.6, add "... for a particular Virtual Connection" at the end of the first sentence for clarification.
9. In 4.4.1, added the Request\_To\_Receive\_Response Operation as a bullet.
10. In 4.4.3, changed the length of the T\_len parameter from 32 bits to 64 bits. Also, rather than use a string of 16 zeros to indicate an unlimited size, said "T\_len = all zeros". Changed the last sentence of the 2nd paragraph to "Note that Out\_of\_Order is necessary for selective retransmission to correct flawed Blocks, otherwise go-back-N retransmission must be used."
11. In 4.4.5, changed "... (in bytes) ..." to "... (the number of bytes in a Block) ..." for clarity. Changed the maximum Blocksize from  $2^{32}$  to  $2^{64}$  bytes.
12. Added all of 4.4.9 detailing the "Opaque data".
13. In 4.4.10, in the last portion deleted the two bullets that said "STU Size  $\leq$  Blocksize" and "STU Size  $\leq$  Bufsize".
14. In table 1, changed the response for a Request\_To\_Receive from "State\_Response" to "Request\_To\_Receive\_Response".
15. Figure 6 was replaced – the original figure had been copied from HIPPI-6400-PH and was not appropriate for -ST. The new figure is HIPPI-ST specific.
16. In 6.1, under T\_len, changed its function in RTS, RTSR, and RTR to just hold the high-order portion of Transfer length. A similar change was made under B\_num, but here B\_num holds the low-order portion of the Transfer length.
17. In 5.1, under OS\_Bufx and OS\_Offset, changed "...ULP parameters" to "...ULP peer-to-peer information (see 4.4.9)".
18. In 6.1, under T\_len, deleted "Request\_To\_Send\_Response" from the list of things that use T\_len.
19. In 6.2, under Send\_State, added the sentence "Send\_State is always valid on Control Operations.". Fleshed out the next sentence so that it clarifies that Data

- Operations need a flag bit to validate Send\_State.
20. In 7.2, under Effect, added the last sentence about how rejection of Concatenate or Source\_Concatenate is indicated.
  21. In 8.1 and 8.3, changed the T\_len parameter in the semantics list from "T\_len" to "T\_len [T\_len, B\_num]". Changed the text for T\_len from "T\_len shall..." to "T\_len, carried in the concatenation of T\_len and B\_num, shall...".
  22. In 8.2, changed the Blocksize parameter from specifying the number of STUs in a Block to the size of the Block as a power of 2. Deleted the T\_len parameter from the Semantics list and the description.
  23. In 8.3, in "Effect", changed "...Reject = 1 in a State\_Response..." to "...Reject = 1 in a Request\_To\_Receive\_Response...".
  24. Added all of 8.4 for the Request\_To\_Receive\_Response Operation. This resulted in all of the following operations increasing their Op code by 1, and their paragraph number by .1.
  25. In 8.5, deleted the T\_len parameter from the semantics list, and from the parameter descriptions. Changed the Blocksize parameter from specifying the number of STUs, to specifying the size.
  26. In 8.6, added a reference for "Opaque data" and deleted the text about it being transferred unmodified.
  27. In 8.8, deleted the sentence about Request\_State\_Response also being the response to a Request\_To\_Receive. Changed "...for additional Operations..." to "...for additional Operations on this Virtual Connection...".
  28. In table 3, changed the T\_len parameter to cover both the T\_len and B\_num fields in RTS and RTR Operations. Deleted the T\_len parameter in RTSR and CTS Operations. Added the RTRR Operation, its abbreviation, and changed the Op code to shoe-horn this new Operation in.
  29. In 9, changed "...owns the Port" to "...is using the Port".
  30. In Table 4, change "State\_Response" to "Request\_To\_Receive\_Response" as one of the responses to Request\_To\_Receive.
  31. In clause 9, added lots of back references for the things being checked.
  32. In 9.4.4, added a check for too large a Bufsize, and added some text to remind people that Bufsize was in powers of 2.
  33. In 9.5.4 and 9.5.5, changed "The Operation shall be discarded..." to "The Operation shall be rejected...".
  34. In 9.5.7, added a new check for Illegal Blocksize, and made it a "should" rather than a "shall".
  35. In 9.5.8, added the CTS Operation needed to set up for an RTR.
  36. In A.2, deleted the last sentence about "The 12 least significant bits of the 48-bit....". Added an open issue about the address mapping between HIPPI-800 and HIPPI-6400 not being resolved.
  37. In figures A.1 and A.2, changed the size of the payload from  $2^{11}$  to  $2^{31}$  bytes, and changed "D\_ULA" to "EtherType".
  38. Added a new revision of Annex C (courtesy of James Hoffman). Since this is all essentially new, margin bars and highlighting were not used.

## Contents

	Page
Foreword (This foreword is not part of American National Standard X3.xxx-199x.) .....	vii
Introduction .....	viii
1 Scope .....	1
2 Normative references.....	1
3 Definitions and conventions .....	2
3.1 Definitions .....	2
3.2 Editorial conventions .....	3
3.2.1 Binary notation .....	3
3.2.2 Hexadecimal notation .....	3
3.3 Acronyms and other abbreviations .....	3
4 System overview .....	3
4.1 Control Channels and Data Channels .....	3
4.2 System model .....	4
4.3 Virtual Connections .....	6
4.3.1 Sequences and Operations.....	6
4.3.2 Ports.....	6
4.3.3 Keys .....	7
4.3.4 Buffer size (Bufsize) .....	7
4.3.5 Max-STU size .....	7
4.3.6 Slots and Sync parameter.....	7
4.3.7 Concatenate .....	8
4.3.8 Source_Concatenate .....	8
4.3.9 Persistent .....	8
4.4 Data movement.....	9
4.4.1 Sequences and Operations.....	9
4.4.2 Transfer identifiers (R_id and S_id) .....	9
4.4.3 Transfer length (T_len) .....	9
4.4.4 Blocks.....	9
4.4.5 Blocksize .....	10
4.4.6 STUs .....	10
4.4.7 Bufx and Offset .....	10
4.4.8 OS_Bufx and OS_Offset .....	10
4.4.9 Opaque data.....	10
4.4.10 Packing examples .....	11
4.5 Operations management.....	12
4.5.1 Flow control .....	12
4.5.2 Status Operations.....	12
4.5.3 Rejected Operations.....	12
4.5.4 Lost Operations .....	12
4.5.5 Interrupts .....	12
5 Service interface.....	13
5.1 Service primitives.....	13
5.2 Sequences of primitives .....	13
6 Schedule Header .....	14
6.1 Schedule Header fields .....	14
6.2 Scheduled Transfer flags.....	15

7	Virtual Connection management .....	16
7.1	Request_Port .....	16
7.2	Request_Port_Response .....	17
7.3	Port_Teardown .....	17
7.4	Port_Teardown_ACK .....	18
7.5	Port_Teardown_Complete .....	18
8	Data movement .....	19
8.1	Request_To_Send .....	19
8.2	Request_To_Send_Response .....	20
8.3	Request_To_Receive .....	20
8.4	Request_To_Receive_Response .....	21
8.5	Clear_To_Send .....	21
8.6	Data .....	22
8.7	Request_State .....	23
8.8	Request_State_Response .....	23
8.9	END .....	24
8.10	END_ACK .....	25
9	Error processing .....	26
9.1	Operation timeout .....	27
9.2	Operation Pairs .....	27
9.3	Syntax errors .....	27
9.3.1	Undefined Opcode .....	27
9.3.2	Unexpected Opcode .....	27
9.4	Virtual Connection errors .....	27
9.4.1	Invalid Key or Port .....	27
9.4.2	Slots exceeded .....	28
9.4.3	Unknown EtherType .....	28
9.4.4	Illegal Bufsize .....	28
9.4.5	Illegal STU size .....	28
9.5	Scheduled Transfer errors .....	28
9.5.1	Invalid S_id .....	28
9.5.2	Bad Data Channel specification .....	28
9.5.3	Concatenate not available .....	28
9.5.4	Source_Concatenate not available .....	28
9.5.5	Persistent not available .....	28
9.5.6	Out of Range B_num, Bufx, Offset, S_count, or Sync .....	29
9.5.7	Illegal Blocksize .....	29
9.5.8	Request_To_Receive problem .....	29
9.5.9	Undefined Flag .....	29

## Tables

Table 1 – Response to a rejected Operation .....	12
Table 2 – Virtual Connection Operations summary between end devices A and B .....	25
Table 3 – Data transfer and status Operations summary between end devices S and R .....	26
Table 4 – Operation pairs guarded by Op_timeout .....	27
Table 5 – Summary of logged errors .....	29

## Figures

Figure 1 – System overview .....	4
Figure 2 – HIPPI-ST over different media .....	4
Figure 3 – Information hierarchy .....	4
Figure 4 – Scheduled Transfer Final Destination model .....	5
Figure 5 – Data packing examples .....	11
Figure 6 – HIPPI-ST service interface .....	13
Figure 7 – Schedule Header contents .....	14
Figure 8 – Flags summary .....	15
Figure A.1 – HIPPI-ST Operations carried in HIPPI-6400-PH Messages .....	31
Figure A.2 – HIPPI-ST Operations carried in HIPPI-FP packets .....	32
Figure B.1 – Many-to-one striping .....	34
Figure B.2 – One-to-many striping .....	34
Figure B.3 – Many-to-many striping .....	34

## Annexes

A Using lower layer protocols .....	30
A.1 HIPPI-6400-PH as the lower layer .....	30
A.2 HIPPI-FP as the lower layer .....	30
B HIPPI-ST striping .....	33
B.1 Striping principles .....	33
B.2 Many-to-one striping .....	33
B.3 One-to-many striping .....	33
B.4 Many-to-many striping .....	34
C Scheduled Transfer example .....	35
C.1 Detailed simple transfer example .....	35
C.1.1 Virtual Connection set up .....	35
C.1.2 Scheduled Transfer set up .....	36
C.1.3 Block 0 transfer .....	36
C.1.4 Block 1 Clear_To_Send .....	37
C.1.5 Ending the Virtual Connection .....	38
C.2 Persistent memory example .....	40
C.2.1 Set up .....	40
C.2.2 Reading .....	41
C.2.3 Writing .....	41
C.2.4 Closing the persistent memory .....	41
C.3 Translated, striped Ethernet to HIPPI-6400 example .....	42
C.3.1 Virtual Connection Setup .....	42
C.3.2 Sending to Hobbs .....	43
C.3.3 Sending to Calvin .....	44

**Foreword** (This foreword is not part of American National Standard X3.xxx-199x.)

This American National Standard specifies a data transfer protocol that uses small control messages to pre-arrange data movement. Buffers are allocated at each end before the data transmission, allowing full-rate, non-congesting data flow between the end devices. The control and data may use different physical media, or may share a single physical medium. Procedures are provided for moving data over HIPPI and other media.

This standard provides an upward growth path for legacy HIPPI-based systems.

This document includes annexes which are informative and are not considered part of the standard.

Requests for interpretation, suggestions for improvement or addenda, or defect reports are welcome. They should be sent to the X3 Secretariat, Information Technology Industry Council, 1250 Eye Street, NW, Suite 200, Washington, DC 20005.

This standard was processed and approved for submittal to ANSI by Accredited Standards Committee on Information Processing Systems, X3. Committee approval of the standard does not necessarily imply that all committee members voted for approval. At the time it approved this standard, the X3 Committee had the following members:

(List of X3 Committee members to be included in the published standard by the ANSI Editor.)

Subcommittee X3T11 on Device Level Interfaces, which developed this standard, had the following participants:

(List of X3T11 Committee members, and other active participants, at the time the document is forwarded for public review, will be included by the Technical Editor.)

## **Introduction**

This American National Standard specifies a data transfer protocol that uses small control messages to pre-arrange data movement. Buffers are allocated at each end before the data transmission, allowing full-rate, non-congesting data flow between the end devices. The control and data may use different physical media, or may share a single physical medium. Procedures are provided for moving data over HIPPI and other media.

Characteristics of a HIPPI-ST include:

- A hierarchy of data units (Scheduled Transfer Units (STUs), Blocks, and Transfers).
- Support for Get and Put Operations.
- Parameters exchanged between end devices for port selection, transfer identification, and Operation validation.
- Features supporting efficient mapping between the sender's and receiver's natural buffer sizes.
- Provisions for resending partial Transfers for error recovery.
- Mappings onto HIPPI-6400-PH, HIPPI-FP (for HIPPI-800 traffic), and Ethernet lower-layer protocols.
- Mappings from IPv4, IPv6, and MPI upper-layer protocols onto Scheduled Transfer.



## American National Standard for Information Technology –

# High-Performance Parallel Interface – Scheduled Transfer (HIPPI-ST)

## 1 Scope

This American National Standard specifies a data transfer protocol that uses small control messages to pre-arrange data movement. Buffers are allocated at each end before the data transmission, allowing full-rate, non-congesting data flow between the end devices. The control and data may use different physical media, or may share a single physical medium. Procedures are provided for moving data over HIPPI and other media.

Specifications are included for:

- Virtual Connection setup and teardown;
- determining the number of Operations the other end can accept;
- determining the buffer size of the other end;
- exchanging Key, Port, transfer identifiers, and buffer size values specific to the end nodes;
- determining a maximum size transmission unit that will not overrun receiver buffer boundaries;
- using buffer indices and 64-bit addresses;
- acknowledging partial transfers so that buffers can be reused;
- providing means for resending partial Transfers for error recovery; and
- terminating transfers in progress.

Note that parts of the Scheduled Transfer protocol depend upon in-order delivery by the lower layer, which may not be available on all media.

## 2 Normative references

The following American National Standards contain provisions which, through reference in this text, constitute provisions of this American National Standard. At the time of publication, the editions indicated were valid. All standards are subject to revision, and parties to agreements based on this standard are encouraged to investigate the possibility of applying the most recent editions of the standards listed below.

ANSI X3.183-1991, *High-Performance Parallel Interface – Mechanical, Electrical, and Signalling Protocol Specification (HIPPI-PH)*

ANSI X3.210-1992, *High-Performance Parallel Interface – Framing Protocol (HIPPI-FP)*

ANSI X3.xxx-199x, *High-Performance Parallel Interface – 6400 Mbit/s Physical Layer (HIPPI-6400-PH)*

ANSI/IEEE Std 802-1990, *IEEE Standards for Local and Metropolitan Area Networks: Overview and architecture (formerly known as IEEE Std 802.1A, Project 802: Local and Metropolitan Area Network Standard — Overview and Architecture).*

ISO/IEC 8802-2:1989 (ANSI/IEEE Std 802.2-1989), *Information Processing Systems – Local Area Networks – Part 2: Logical link control.*

### 3 Definitions and conventions

#### 3.1 Definitions

For the purposes of this standard, the following definitions apply.

**3.1.1 Block:** An ordered set of one or more STUs within a Scheduled Transfer. (See figure 3 and 4.4.4.)

**3.1.2 Buffer Index (Bufx):** A 32-bit parameter identifying the starting address of a data buffer. Bux may be either a pointer to the starting address or the most significant part of a 64-bit starting address.

**3.1.3 Concatenate:** An addressing mode using 64-bit addresses rather than buffer indices. (See 4.3.7.)

**3.1.4 Control Channel:** The logical channel that carries the Control Operations.

**3.1.5 Control Operation:** A control function consisting of a Schedule Header and an optional 32-byte payload. (See figure 3.)

**3.1.6 Data Channel:** The logical channel that carries the data payload.

**3.1.7 Data Operation:** A data movement Operation consisting of a Schedule Header and up to 2 gigabytes of user payload. (See figure 3.)

**3.1.8 Final Destination:** The end device that receives, and operates on, the data payload. This is typically a host computer system, but may also be a non-transparent translator, bridge, or router.

**3.1.9 Key:** A local identifier used to validate Operations. (See 4.3.3.)

**3.1.10 log:** The act of making a record of an event for later use.

**3.1.11 Opaque data:** Eight bytes of Source ULP to Destination ULP peer-to-peer information carried separately from the data payload. (See 4.4.9)

**3.1.12 Operation:** A Scheduled Transfer function, i.e., a Control Operation or the data movement specified in an STU.

**3.1.13 optional:** Characteristics that are not required by HIPPI-ST. However, if any optional characteristic is implemented, it shall be

implemented as defined in HIPPI-ST.

**3.1.14 Originating Source:** The end device that generates the data payload. This is typically a host computer system, but may also be a non-transparent translator, bridge, or router.

**3.1.15 Persistent:** A control mode used to retain buffers for multiple Transfers. (See 4.3.9.)

**3.1.16 Port:** A logical connection within an end device. (See 4.3.2.)

**3.1.17 Scheduled Transfer:** An information transfer, normally used for bulk data movement and low processing overhead, where the Originating Source and Final Destination prearrange the transfer using the protocol defined in this standard.

**3.1.18 Scheduled Transfer Unit (STU):** The data payload portion of a Data Operation moved from an Originating Source to a Final Destination. STUs are the basic components of Blocks. (See figure 3 and 4.4.6.)

**3.1.19 Slot:** A space reserved for a Control Operation, or the Schedule Header portion of an STU, in the end device. (See 4.3.6.)

**3.1.20 Transfer:** An ordered set of one or more Blocks within a Scheduled Transfer. (See figure 3 and 4.2.)

**3.1.21 upper-layer protocol (ULP):** The protocol above the service interface. These could be implemented in hardware, software, or they could be distributed between the two.

**3.1.22 Virtual Connection:** A bi-directional logical connection used for Scheduled Transfers between two end devices. A Virtual Connection contains a logical Control Channel and a logical Data Channel in each direction.

### 3.2 Editorial conventions

In this standard, certain terms that are proper names of signals or similar terms are printed in uppercase to avoid possible confusion with other uses of the same words (e.g., END). Any lowercase uses of these words have the normal technical English meaning.

A number of conditions, sequence parameters, events, states, or similar terms are printed with the first letter of each word in uppercase and the rest lowercase (e.g., Block, Transfer). Any lowercase uses of these words have the normal technical English meaning.

The word *shall* when used in this American National standard, states a mandatory rule or requirement. The word *should* when used in this standard, states a recommendation.

#### 3.2.1 Binary notation

Binary notation is used to represent relatively short fields. For example a two-bit field containing the binary value of 10 is shown in binary format as b'10'.

#### 3.2.2 Hexadecimal notation

Hexadecimal notation is used to represent some fields. For example a two-byte field containing a binary value of b'11000100 00000011' is shown in hexadecimal format as x'C403'.

### 3.3 Acronyms and other abbreviations

<b>ACK</b>	acknowledge indication
<b>CTS</b>	Clear_To_Send
<b>END_A</b>	END_ACK
<b>HIPPI</b>	High-Performance Parallel Interface
<b>K</b>	kilo ( $2^{10}$ or 1024)
<b>KB</b>	kilobyte (1024 bytes)
<b>MAC</b>	Media Access Control
<b>PT</b>	Port_Teardown
<b>PTA</b>	Port_Teardown_ACK
<b>PTC</b>	Port_Teardown_Complete
<b>RQP</b>	Request_Port
<b>RQPR</b>	Request_Port_Response
<b>RS</b>	Request_State
<b>RSR</b>	Request_State_Response
<b>RTR</b>	Request_To_Receive

<b>RTRR</b>	Request_To_Receive_Response
<b>RTS</b>	Request_To_Send
<b>RTSR</b>	Request_To_Send_Response
<b>STU</b>	Scheduled Transfer Unit
<b>ULP</b>	upper-layer protocol

## 4 System overview

This clause provides an overview of the structure, concepts, and mechanisms used in Scheduled Transfers. Figure 1 gives an example of Scheduled Transfers being used to communicate between device A and device B over some physical media. Annex C describes the steps in a typical Scheduled Transfer. Figure 2 shows HIPPI-ST being used over different media.

### 4.1 Control Channels and Data Channels

Each Transfer has an Originating Source and Final Destination. Each Originating Source and Final Destination shall have a Control Channel and one or more Data Channels. The Originating Source sends the payload data, and the Final Destination receives the payload data.

Control Operations shall be exchanged over the Control Channel. Scheduled Transfer Units (STUs), i.e., data payload, shall be exchanged over the Data Channel(s). The information volume on the Data Channel(s) will be probably many times the volume on the Control Channel; hence the available bandwidths should be balanced accordingly. For best performance, the Control Channel should have low latency.

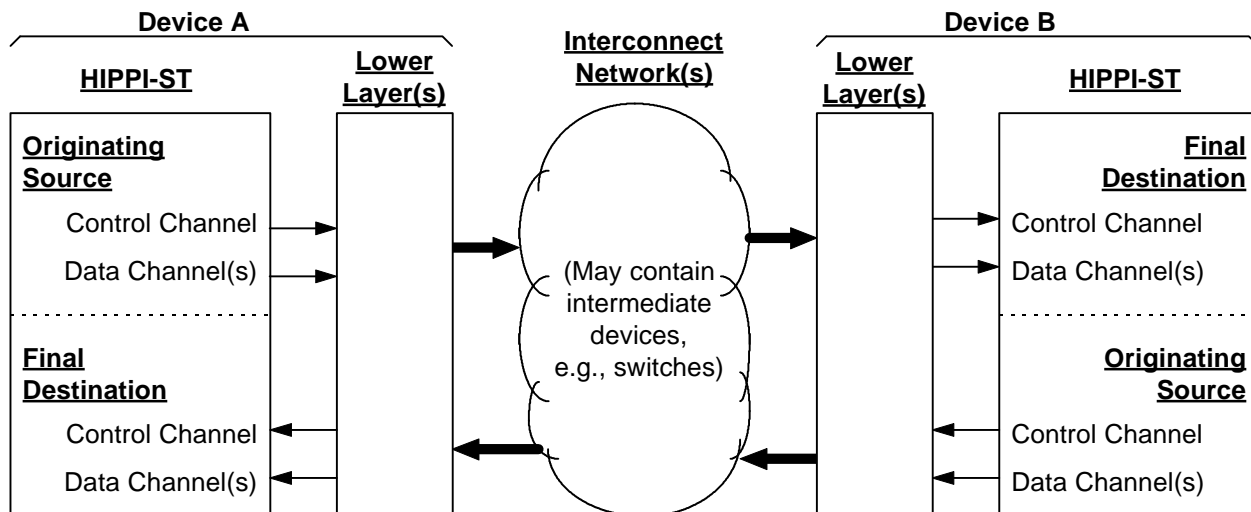


Figure 1 – System overview

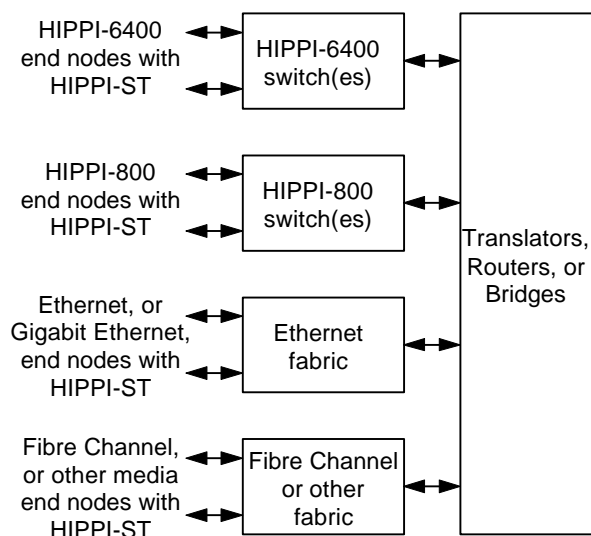
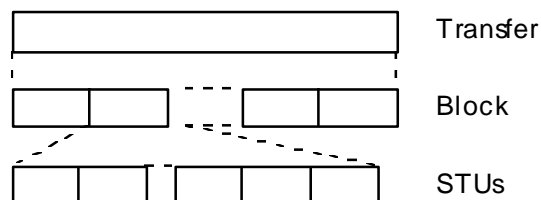


Figure 2 – HIPPI-ST over different media

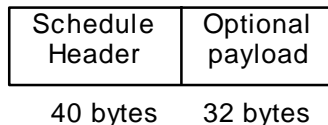
#### 4.2 System model

Multiple write (Put) or read (Get) functions may be executed to move data units, called Transfers, over a Virtual Connection. As shown in figure 3, a Transfer is composed of one or more Blocks, and Blocks are composed of one or more STUs. The Scheduled Transfer protocol shall package the Transfer in Blocks and STUs for delivery using lower layer protocol(s) and media. An STU

shall be the data payload portion of a Data Operation. A Data Operation shall consist of a 40-byte Schedule Header and an STU of up to 2 gigabytes ( $2^{31}$  bytes). A Control Operation shall consist of a 40-byte Schedule Header, and may contain an additional 32 bytes of optional payload.



#### Control Operation



#### Data Operation

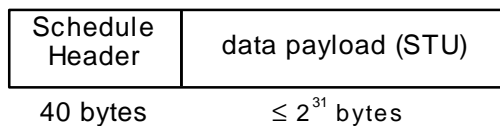


Figure 3 – Information hierarchy

Figure 4 shows the model used on a Final Destination for the Scheduled Transfers. The model on an Originating Source would be similar.

As Control Operations and Data Operations are received, the Schedule Header of each is placed in the Schedule Header queue for execution. State information about the number of empty Slots in the queue is available to the other end so that it can avoid overrunning the queue.

The Virtual Connection Descriptor contains:

- static parameters defining the Virtual Connection from the view of both the remote end device and local end device (the top portion of the Virtual Connection Descriptor box in figure 4);
- current state information about the number of empty "Slots" for Operation Schedule Headers, and Operation Retry and Timeout parameters;
- identifiers for each of the Virtual Connection's Transfers.

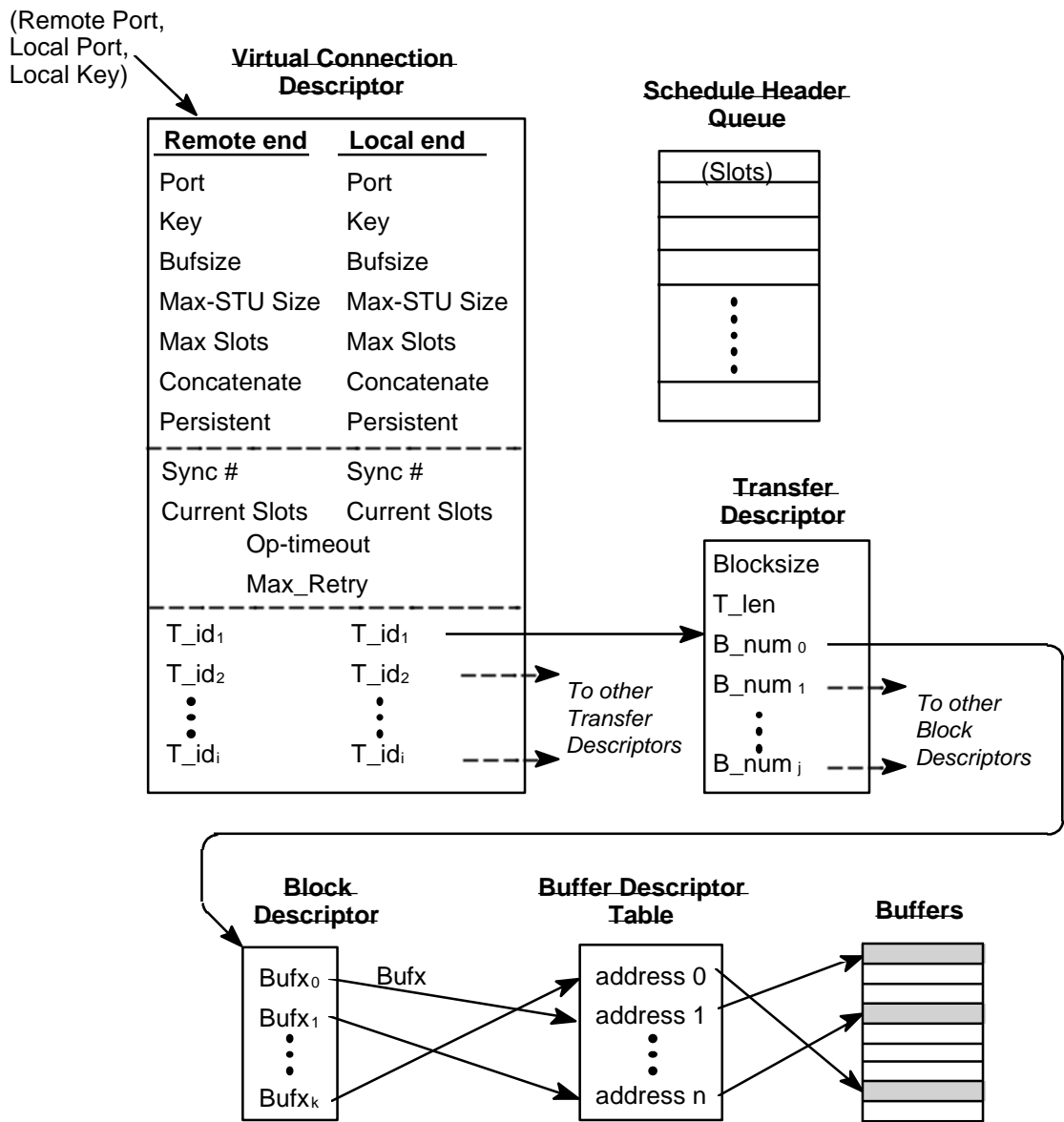


Figure 4 – Scheduled Transfer Final Destination model

A Transfer Descriptor, for each Transfer, contains the Transfer length (T\_len, in bytes), the Blocksize (in bytes), and includes pointers to Block Descriptors. The Block Descriptors (one for each Block of a Transfer) identify the set of contiguous Buffer Index (Bufx) values assigned to the Block. And finally, the Buffer Descriptor Table provides a base memory address for each Bufx.

In an effort to achieve maximum transfer rates and efficiency, the receiver's job is made as easy as possible, even at the expense of the transmit side. It is expected that after validating an Operation in the Final Destination, only a single lookup will be needed to derive the absolute memory address and correctly place the data.

### 4.3 Virtual Connections

Scheduled Transfers between an Originating Source and Final Destination are pre-arranged to decrease computational overhead during the Transfer by allocating buffers at each end device. The bi-directional path between the end devices is called a Virtual Connection. A Virtual Connection shall consist of an Originating Source and Final Destination in each end device.

Once the Final Destination has indicated its ability to accept the STUs, the Virtual Connection should not become congested. In essence, the Final Destination smoothly controls the flow. For comparison, without pre-arranging the buffers, the Originating Source would blindly send data into the interconnection network where it might have to wait for buffers to be assigned in the Final Destination. On the down-side, Scheduled Transfers require additional Control Operations and round-trip latency. Once established, a Virtual Connection may be used to carry multiple Transfers. This Scheduled Transfer protocol does not handle network resource reservations.

#### 4.3.1 Sequences and Operations

During Virtual Connection setup, the end devices shall exchange parameters specific to each device. These parameters, shown in the upper portion of the Virtual Connection Descriptor box in figure 4 and detailed below, include values for:

- Port numbers (e.g., a Port dedicated to HIPPI-FP or IP traffic);
- Keys (used for authenticating Operations);
- native buffer sizes (Bufsize) for determining Final Destination buffer tiling;
- maximum STU size;
- maximum number of outstanding Operations (Slots) to keep from overflowing the command queues;
- whether or not they support Concatenate and Persistent modes.

The parameters assigned during setup shall apply for the life of the Virtual Connection. Once established, the Virtual Connection is accessed as shown in figure 4 by the tuple "remote Port", "local Port", and "local Key". The Control Operations defined for Virtual Connection setup are:

- Request\_Port (See 7.1.)
- Request\_Port\_Response (See 7.2.)

The Control Operations defined for Virtual Connection teardown are:

- Port\_Teardown (See 7.3.)
- Port\_Teardown\_ACK (See 7.4.)
- Port\_Teardown\_Complete (See 7.5.)

#### 4.3.2 Ports

Ports identify upper-layer entities within an end device. The Port values shall be assigned by the local end device and have no meaning on the other end device. For example, when end device A requests a Virtual Connection to end device B, A shall select the value for A-Port and shall send it to B in the Request\_Port Operation. B shall store the A-Port value and shall return it to A in every Operation over this Virtual Connection. Likewise, B shall select the value for B-Port.

An exception is the "well-known Port", i.e., Port x'0000'. In this case, a request sent to the "well-known Port" shall result in the receiving end device assigning a specific local Port value based on the EtherType parameter. EtherType parameter values shall be as assigned in the current "Assigned Numbers" RFC, e.g., RFC 1700<sup>1)</sup>. For example, if the HIPPI-ST is used to

encapsulate TCP/IP, then the EtherType would be x'0800'. If HIPPI-ST is being used to encapsulate legacy HIPPI-FP or user data, then the EtherTypes would be x'8180' and x'8181' respectively.

If the incoming Port number is invalid, then the Operations shall not be executed (see 9.4.1). A Port value of x'0000' is valid in Request\_Port Operations; invalid in all other Operations.

### 4.3.3 Keys

Like the Ports, each end device shall select its own 32-bit Key value for use on the Virtual Connection. For example, when end device A requests a Virtual Connection to end device B, A shall select the value for A-Key and shall send it to B in the Request\_Port Operation. B shall store the A-Key value and shall return it to A in every Operation over this Virtual Connection. The A-Key value has no meaning in B; it is only significant in A where it shall be used to validate that the Operation presented is really associated with this Virtual Connection. Likewise, B shall select the value for B-Key. Keys are similar in nature to passwords; if the Key doesn't match, then the Operation shall not be executed (see 9.4.1).

### 4.3.4 Buffer size (Bufsize)

Each end shall define the buffer size, in bytes, that it wants to use. Buffer sizes may be the same as host page sizes. It is most efficient when the buffer sizes are the same on both ends, but differing buffer sizes are supported (see annex C). The buffer sizes shall be  $\geq 256$  bytes and shall be an integral power of two, i.e.,  $2^x$  where  $8 \leq x \leq 64$ .

### 4.3.5 Max-STU size

The Max-STU size, exchanged during Virtual Connection setup, establishes the maximum data payload size of an STU (see 4.4.10). Each end device declares the desired Max-STU size it is prepared to receive. The Max-STU size must be

no larger than its Bufsize. Intermediate devices with smaller buffer sizes may lower this value. Note that the Max-STU size in each direction may be different.

Additionally, the Max-STU size shall be  $\geq 256$  bytes and an integral power of two i.e.,  $2^x$  where  $8 \leq x \leq 32$ .

### 4.3.6 Slots and Sync parameter

The term Slot denotes memory at an end device reserved for storing the Schedule Headers of incoming Operations for a particular Virtual Connection. Each Operation arriving at an end device consumes one Slot, except for Data Operations which consume a Slot only if a Silent = 0 or Interrupt = 1. An Originating Source shall control the flow of Operations by sending no more Operations than there are Slots available at the other end. Any Operations that are sent in excess of the number of available Slots may be discarded by the receiver (see 9.4.2). The end device supplying the Slot information shall advertise at least one too few Slots, i.e., keeping one in reserve for an END or Request\_State\_Response Operation to prevent deadlocks.

An end device learns the initial number of Slots available (Slots value) at the remote end device during the Virtual Connection setup (see 7.1 and 7.2). Later, an end device obtains the current Slots value by reading the Slots parameter in a received Request\_State\_Response. An end device may solicit a Request\_State\_Response from the remote end by either of two methods: by setting the Send\_State flag in the Schedule Header of a Data Operation, or by sending a Request\_State Operation. A received Slots value of x'FFFFFFFF' indicates that the remote end does not implement Slot accounting.

NOTE – Slot flow control may not be needed when the maximum number of Control Operations is otherwise bounded or where dropped Operations are acceptable.

<sup>1)</sup> RFC (Request For Comment) documents are working standards documents from the TCP/IP internetworking community. Copies of these documents are available from numerous electronic sources (e.g., <http://www.ietf.org>) or by writing to IETF Secretariat, c/o Corporation for National Research Initiatives, 1895 Preston White Drive, Suite 100 Reston, VA 20191-5434, USA.

The received Slot value is a snapshot of the number of Slots available at the remote end device when the remote end device received the soliciting Operation. The local end device may continue to send Operations after soliciting a **Request\_State\_Response** and may also solicit multiple responses before receiving a reply. The lower bound on the number of available Slots at the remote end device is determined by the local end device which adjusts its vision of the number of Slots to account for outstanding Operations. The adjustment consists of subtracting, from the number of Slots indicated in the received **Request\_State\_Response** Operation, the number of Slot-consuming Operations sent by the local end device after a **Request\_State\_Response** solicitation.

The local end device can use the Sync parameter to identify **Request\_State\_Response** messages when there are multiple outstanding solicitations. The Sync parameter in a Data or Request\_State shall be copied and returned by the remote end device in the corresponding **Request\_State\_Response**. The Sync parameter may be used by the local end device to mark the request, and thus identify the **Request\_State\_Response** with a particular solicitation. The Sync values are locally determined.

#### 4.3.7 Concatenate

The Concatenate flag (see 6.2) controls the addressing mode for the Bufx and Offset parameters.

- When Concatenate = 0, Bufx shall specify a Buffer Index for placing the data in the Final Destination.
- When Concatenate = 1, the Bufx and Offset fields shall be concatenated into a single 64-bit address, with Bufx containing the most-significant bytes of the address (see 4.4.7).

The value of the Concatenate flag shall be consistent for an entire Transfer, i.e., switching back and forth between 64-bit addresses and Buffer Indexes within a Transfer is not allowed (see 4.4.7).

Concatenate is only usable between hosts that mutually agree. Agreement is reached by

controlling the Concatenate flag bit during the Virtual Connection setup (see 7.1 and 7.2).

#### 4.3.8 Source\_Concatenate

The Source\_Concatenate flag (see 6.2) controls the addressing mode for the OS\_Bufx (Originating Source Buffer index) and OS\_Offset parameters. Note that OS\_Bufx and OS\_Offset control the data addressing on the Originating Source and are only used in the Request\_To\_Receive Operation.

- When Source\_Concatenate = 0, OS\_Bufx shall specify a Buffer Index pointing to the data in the Originating Source (see 4.4.8).
- When Source\_Concatenate = 1, the OS\_Bufx and OS\_Offset fields shall be concatenated into a single 64-bit address, with OS\_Bufx containing the most-significant bytes of the address (see 4.4.8).

Source\_Concatenate is only usable between hosts that mutually agree. Agreement is reached by controlling the Source\_Concatenate flag bit during the Virtual Connection setup (see 7.1 and 7.2).

#### 4.3.9 Persistent

The Persistent flag (see 6.2) controls buffer retention in the Final Destination for the Virtual Connection.

- When Persistent = 1, the memory in the Final Destination allocated for the Scheduled Transfer shall be retained for multiple transfers and not released until a Port\_Teardown or an END Operation occurs. Note that Persistent = 1 bypasses the flow control provided by Clear\_To\_Send, i.e., a Data Operation may be sent at any time whether or not a Clear\_To\_Send Operation has been received. Sending information to a Frame Buffer is an example of where Persistent might be used.
- When Persistent = 0, the memory for a Block may be allocated for other uses after the Block is complete. All Data Operations must be enabled by a Clear\_To\_Send or Request\_To\_Receive Operation.

Persistent is only usable between hosts that mutually agree. Agreement is reached by controlling the Persistent flag bit during the



Virtual Connection setup (see 7.1 and 7.2).

## 4.4 Data movement

### 4.4.1 Sequences and Operations

A write data sequence (which may be initiated by either end of the Virtual Connection) shall be setup by the end devices exchanging transfer identifiers (T\_id's), specific to each device, and length parameters. The Control Operations setting up a write data sequence are:

- Request\_To\_Send (See 8.1.)
- RTS\_Response (See 8.2.)

A read data sequence, which moves the Transfer as a single Block, requires that both ends had previously allocated resources for the entire read sequence with a Request\_To\_Send. The Control Operations setting up a read data sequence are:

- Request\_To\_Receive (See 8.3.)
- Request\_To\_Receive\_Response (See 8.4.)

The Final Destination controls the data flow with:

- Clear\_To\_Send (See 8.5.)

Data payloads for the read and write data movements are carried in STUs. STUs are sent with:

- Data (See 8.6.)

State information can be requested in a Data Operation or with a Request\_State Control Operation.

- Request\_State (See 8.7.)
- Request\_State\_Response (See 8.8.)

The Control Operations below are used to abort limited size Transfers. Unlimited size Transfers shall use this method to signal the end of the Transfer.

- END (See 8.9.)
- END\_ACK (See 8.10.)

### 4.4.2 Transfer identifiers (R\_id and S\_id)

Like the Ports and Keys, each end device shall also select its own non-zero 16-bit Transfer

identifier (T\_id) value for a data movement on the Virtual Connection. For example, when end device S requests to write to end device R, S shall select the value for its T\_id and shall send it to R in the Request\_To\_Send Operation. R shall store S's T\_id value and shall return it to S in every Operation concerning this Transfer. Likewise, R shall select its T\_id value and send it to S in a Request\_To\_Send\_Response or Clear\_To\_Send Operation. For each Operation, the sender shall put its T\_id in the S\_id field and the receiver's T\_id value in the R\_id field.

NOTE – The Virtual Connection is symmetrical; either end device may initiate a data movement. For example, S could be end device A that initiated the Virtual Connection setup, or it could be end device B. Different names were used for clarity.

### 4.4.3 Transfer length (T\_len)

The 64-bit Transfer length parameter (T\_len) specifies the total number of data payload bytes in the Transfer. T\_len does not include the Schedule Header or any lower-layer headers. T\_len = all zeros shall indicate an unlimited size Transfer. An unlimited size Transfer is terminated by an END Operation (see 8.9).

### 4.4.4 Blocks

Scheduled Transfer flow control, striping, acknowledgments, and resource allocation are all done on a Block basis. Block numbers (B\_num) shall be numbered starting at zero and shall increment by one for each following Block.

Blocks comprising a Transfer shall be enabled for transmission in sequential order unless both the Originating Source and Final Destination indicated Out\_of\_Order capability during the Virtual Connection setup. Note that Out\_of\_Order is necessary for selective retransmission to correct flawed Blocks, otherwise go-back-N retransmission must be used.

Request\_State\_Response Operations indicate the highest numbered Block received correctly by the Final Destination. Request\_State\_Response Operations can be requested by setting the Send\_State flag bit in Data Operations or by sending Request\_State Operations. In addition, Request\_State Operations can ask if a particular

Block was received correctly. Use of these mechanisms allows the Originating Source to verify correct reception and to identify flawed Blocks for potential retransmission.

#### 4.4.5 Blocksize

The Blocksize (the number of bytes in a Block) for a Transfer is established when a Transfer is initiated, i.e., with a Request\_To\_Send\_Response or Clear\_To\_Send Operation (see 8.2 and 8.5). Blocksize is expressed as a power of two, i.e.,  $2^x$  where  $8 \leq x \leq 64$ . All of the Blocks of a Transfer shall be full size, except for the first and/or last Block of a Transfer which can be smaller (the first Block will be smaller by the initial Offset value, and the last Block will be whatever completes the Transfer).

#### 4.4.6 STUs

The STUs of a Block shall be transmitted in order. STU numbers (S\_count) shall start with zero and increment by one for each following STU. The last STU of a Block shall be marked with Last = 1. No STU shall extend past a Final Destination's buffer boundary, Blocksize boundary, or Transfer boundary.

#### 4.4.7 Bufx and Offset

Bufx contains either a Buffer Index, or the high-order portion of a 64-bit address, in the Final Destination. Selection between these two is controlled by the Concatenate flag (see 4.3.7 and 6.2). If more than one Buffer Index is required for a Block, i.e., buffer size (Bufsize) is less than Blocksize, then the Bufx parameter in the Clear\_To\_Send Operation shall specify the initial Bufx, and any additional Bufx values shall be sequential.

Offset may be used to start at other than the first byte of a Final Destination's buffer. For the first STU of a Block, the Offset value shall be the same as received in the Clear\_To\_Send for the Block. Subsequent STUs of the Block shall adjust the Bufx and Offset values based on the Final Destination's buffer size and the STU size used by the Originating Source.

The Offset value associated with the first block of a Transfer (I\_Offset) is included in all

Clear\_To\_Send Operations. This allows the Originating Source to compute the starting address for any Block without having received the Clear\_To\_Send for the first Block. Clear\_To\_Send Operations can occur out of order, e.g., as the result of striping.

Best performance will usually be achieved when an Offset value of zero is specified. Use of non-zero offset values may degrade performance, depending upon underlying hardware transfer mechanisms.

#### 4.4.8 OS\_Bufx and OS\_Offset

OS\_Bufx specifies either a Buffer Index, or the high-order portion of a 64-bit address, in the Originating Source during a Request\_To\_Receive Operation. Selection between these two is controlled by the Source\_Concatenate flag (see 4.3.8 and 6.2). If more than one Buffer Index is required for a Block, i.e., buffer size < Blocksize, then the OS\_Bufx parameter shall specify the initial Bufx, and any additional Bufx values shall be sequential.

OS\_Offset may be used to start at other than the first byte of a Source buffer. Note that OS\_Bufx and OS\_Offset are only used with Request\_To\_Receive Operations, and Request\_To\_Receive Operations only specify one Block.

#### 4.4.9 Opaque data

Opaque data is eight bytes of ULP peer-to-peer information carried in Data Operations. One Opaque data value shall be used with a Block, i.e., all of the STUs of a particular Block whose header is delivered to the Final Destination's ULP (i.e., Silent = 0, see 6.2) shall use the same Opaque data value. In STUs with Silent = 1, the Opaque data may be any value. Note that the individual Blocks of a Transfer may have different Opaque data values.

A Data Operation shall carry the Opaque data in the Scheduled Header OS\_Bufx and OS\_Offset fields, with OS\_Bufx containing the most-significant bytes. The Opaque data shall be passed unmodified from the Originating Source to the Final Destination. Note that the Opaque data uses Slot resources while the data payload uses Bufx resources. The Opaque data shall not

be counted in the length, tiling, or Bufx calculations.

#### 4.4.10 Packing examples

Figure 5 shows three possibilities for packing the same Transfer into a receiver's buffers. All three examples show a group of seven of the receiver's buffers on the top line. Each buffer is pointed to by a Bufx, and the data in the first buffer starts at an Offset value. The Transfer is the shaded bar, with transmission going from left to right. The Block boundaries are shown above the shaded bar, and the resulting STU boundaries are shown below the shaded bar.

Example (a), at the top, shows the case where the buffers and Blocks are the same size. Notice that the first Block is smaller than the other Blocks by the Offset value. No Offset is required for the other Blocks. The last Block of the Transfer is also smaller, i.e., the Transfer did not end on a Block boundary. While the STU

boundaries lined up nicely, the sender could have used multiple STUs, but the STUs cannot be larger than Max-STU.

Example (b) shows multiple Blocks per receiver buffer. The Blocks that do not start on a buffer boundary would use the Offset parameter to position the data.

Example (c) shows the Blocksize covering two of the receiver's buffers.

In summary, STUs cannot cross Block, buffer, or Transfer boundaries. Relationships include:

$$\text{STU size} \leq \text{Max-STU size}$$

$$\text{Max-STU size} \leq \text{Blocksize}$$

$$\text{Max-STU size} \leq \text{Bufsize}$$

Note that Blocksize can be larger, smaller, or the same as Bufsize.

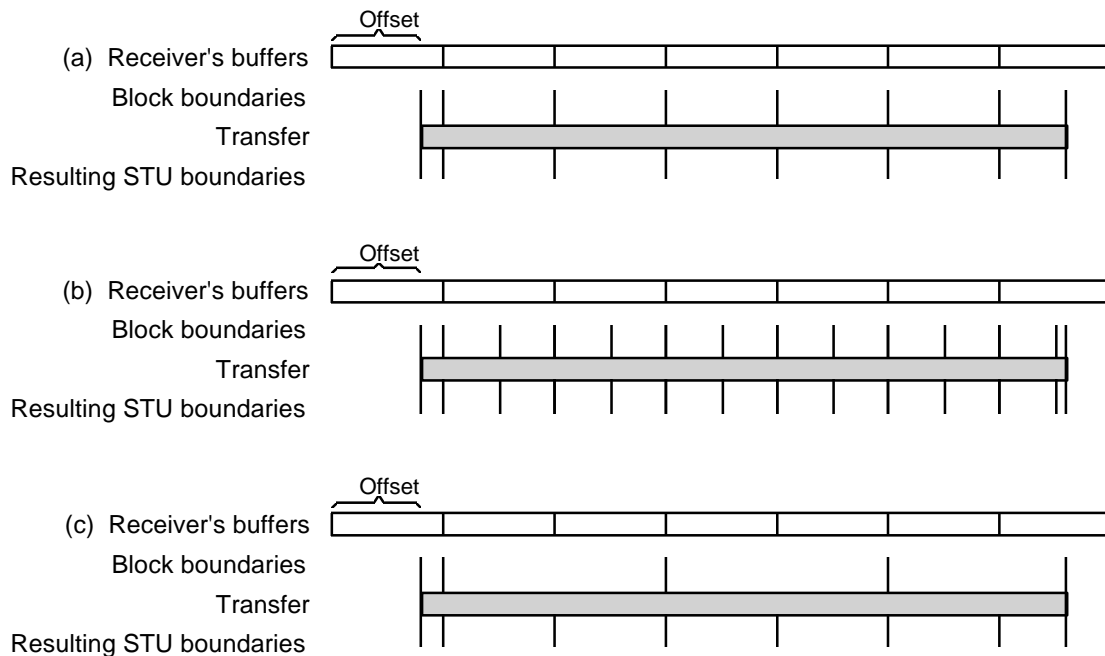


Figure 5 – Data packing examples

## 4.5 Operations management

### 4.5.1 Flow control

Data flow control is achieved with Clear\_To\_Send and Request\_To\_Receive Operations; each one sent by the Final Destination gives the Originating Source permission to send one Block. Flow control is overridden when Persistent = 1; here Data Operations may be sent without having first received Clear\_To\_Send Operations.

Operation flow control is achieved by an Operation's sender not overrunning the Slots value (see 4.3.6).

### 4.5.2 Status Operations

Request\_State (see 8.7) and Request\_State\_Response (see 8.8) Operations are used to request and supply status information about the state of the remote end device. They can be used to see which Blocks have been received correctly and the number of empty Slots available. The Sync parameter (see 4.3.6) is used to provide a common reference point for the local and remote end devices, i.e., to match Request\_State and Request\_State\_Response Operations.

### 4.5.3 Rejected Operations

If the receiving end device is unable to execute an Operation, then the receiving device shall set the Reject flag bit = 1 in the response. Table 1 shows the response when an Operation is rejected. The recovery actions taken when an Operation is rejected are beyond the scope of this standard.

**Table 1 – Response to a rejected Operation**

Rejected Operation	Response (w/ Reject = 1)
Request_Port	Request_Port_Response
Request_To_Send	Request_To_Send_Response
Request_To_Receive	Request_To_Receive_Response

### 4.5.4 Lost Operations

Errors other than syntactic errors are manifested as missing Operations, which occur when the underlying physical medium discards or damages a transmission. Each Scheduled Transfer Operation is defined as part of a two-way handshake or a three-way handshake. Thus, for each command Operation there is a corresponding response Operation, and for some response Operations there is also a corresponding completion Operation.

Each Operation that expects a response is guarded with a timeout whose value is referred to as Op\_timeout (see 9.1). An Operation shall be re-tried up to Max\_Retry times (see 9.1) if the sending end device does not receive the expected response (see 9.2 and table 5).

Data transmissions (i.e., Data Operations) are an exception to this timeout mechanism and are referred to the ULP for resolution (see 9.2).

### 4.5.5 Interrupts

An Interrupt causes a signal to be delivered to the receiving end device ULP. An Interrupt can be requested with any Operation by setting Interrupt = 1.

## 5 Service interface

This clause specifies the services provided by HIPPI-ST. The intent is to allow ULPs to operate correctly with this HIPPI-ST. How many of the services described herein are chosen for a given implementation is up to that implementor; however, a set of HIPPI-ST services must be supplied sufficient to satisfy the ULP(s) being used. The services as defined herein do not imply any particular implementation or any interface.

Figure 6 shows the relationship of the HIPPI-ST interfaces.

### 5.1 Service primitives

The primitives, in the context of the state transitions in clause 5, are declared required or optional. Additionally, parameters are either required, conditional, or optional. All of the primitives and parameters are considered as required except where explicitly stated otherwise.

HIPPI-ST service primitives are of four types.

- *Request primitives* are issued by a service user to initiate a service provided by the HIPPI-ST. In this standard, a second Request primitive of the same name shall not be issued until the Confirm for the first request is received.
- *Confirm primitives* are issued by the HIPPI-ST to acknowledge a Request.
- *Indicate primitives* are issued by the HIPPI-ST to notify the service user of a local event. This primitive is similar in nature to an unsolicited interrupt. Note that the local event may have been caused by a service Request. In this standard, a second Indicate primitive of the same name shall not be issued until the Response for the first Indicate is received.
- *Response primitives* are issued by a service user to acknowledge an Indicate.

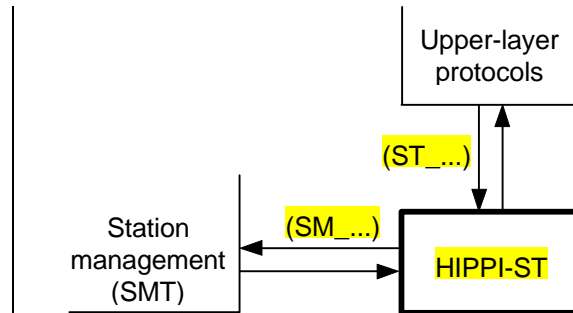


Figure 6 – HIPPI-ST service interface

### 5.2 Sequences of primitives

The order of execution of service primitives is not arbitrary. Logical and time sequence relationships exist for all described service primitives. Time sequence diagrams are used to illustrate a valid sequence. Other valid sequences may exist. The sequence of events between peer users across the user/provider interface is illustrated. In the time sequence diagrams, the HIPPI-ST users are depicted on either side of the vertical bars, while the HIPPI-ST acts as the service provider.

*NOTE - The intent is to flesh out the service primitives similar to what is in HIPPI-PH today.*

#### Service interface considerations -

*(These are notes that have been collected during the document reviews, and should be considered when the service interface is written.)*

*Should there be a priority, or time-to-live, for individual Transfers? On a per connection basis?*

*Pass the full ST header to/from the ULP.*

*Service the slots in order of arrival, i.e., FIFO.*

*Interrupts are passed independent of the Slots, i.e., whenever an Interrupt is put in the Slots queue.*

*There is a Port-basis for the Service Interface*

## 6 Schedule Header

The Schedule Header is shown in figure 7 as a group of 32-bit words. The Schedule Header fields are named for the most common parameter for which the field is used. Many of the fields have different uses depending on the Operation type, and some Operations do not use one or more of the fields at all. The usage for each field is listed below and summarized in tables 2 and 3.

			Bytes
Op	Flags	S_count	00-03
R_Port		S_Port	04-07
Key			08-11
R_id		S_id	12-15
Bufx			16-19
Offset			20-23
T_len			24-27
B_num			28-31
OS_Bufx			32-35
OS_Offset			36-39

Figure 7 – Schedule Header contents

### 6.1 Schedule Header fields

The Schedule Header fields shall be as follows. If an Operation does not use a particular Schedule Header field, then that field shall be transmitted as zeros.

**Op** (5 bits, high-order 5 bits of byte 00) – The Scheduled Transfer Operation. See tables 2 and 3 for a summary of Op values. Unspecified Op values are reserved.

**Flags** (11 bits, low-order 3 bits of byte 00, and all of byte 01) – Control flags (see 6.2).

**S\_count** (16 bits, bytes 02-03):

- In *Request\_Port*, *Request\_Port\_Response*, and *Request\_State\_Response* Operations: the number of available Slots (see 4.3.6);
- In *Request\_To\_Send\_Response* and *Clear\_To\_Send* Operations: the Blocksize parameter (see 4.4.5);

– In *Data Operations*: the STU number (see 4.4.6).

**R\_Port** (16 bits, bytes 04-05) – The receiver's logical Port for this Operation (see 4.3.2).

**S\_Port** (16 bits, bytes 06-07) – The sender's logical Port for this Operation (see 4.3.2).

**Key** (32 bits, bytes 08-11) – Virtual Connection identifier. Generated independently by each end during the Virtual Connection setup. (See 4.3.3.)

**R\_id** (16 bits, bytes 12-13) – The receiver's Transfer identifier for this Operation (see 4.4.2).

**S\_id** (16 bits, bytes 14-15) – The sender's Transfer identifier for this Operation (see 4.4.2).

**Bufx** (32 bits, bytes 16-19):

– In *Request\_Port* and *Request\_Port\_Response* Operations: the maximum buffer size (Bufsize) supported by the end device (see 4.3.4);

– In *Request\_To\_Receive*, *Clear\_To\_Send*, and *Data Operations*: the Buffer Index at the Final Destination or the high-order portion of a 64-bit concatenated address (see 4.4.7).

**Offset** (32 bits, bytes 20-23):

– In *Request\_Port* and *Request\_Port\_Response* Operations: the sender's Key value (see 4.3.3);

– In *Request\_To\_Receive*, *Clear\_To\_Send*, and *Data Operations*: the Final Destination's Offset within a Bufox, or the low-order portion of a 64-bit concatenated address (see 4.4.7);

– In *Request\_State\_Response* Operations: the Block number of the highest numbered contiguous Block received correctly (see 4.4.4).

**T\_len** (32 bits, bytes 24-27):

– In *Request\_Port* and *Request\_Port\_Response* Operations: the Max-STU size (see 4.3.5);

– In *Request\_To\_Send* and *Request\_To\_Receive* Operations: the high-order portion of the length, in bytes, of the Transfer data (see 4.4.3);

– In *Data*, *Request\_State*, and *Request\_State\_Response* Operations: the Sync parameter (see 4.3.6).



**B\_num** (32 bits, bytes 28-31):

– In *Request\_Port Operations*: the EtherType parameter (see 4.3.2);

– In *Request\_To\_Send, Request\_To\_Send Response, and Request\_To\_Receive Operations*: the low-order portion of the length, in bytes, of the Transfer data (see 4.4.3);

– In *Clear\_To\_Send and Data Operations*: the Block number being requested or transmitted (see 4.4.4);

– In *Request\_State and Request\_State\_Response Operations*: the Block number being queried or responded to (see 4.4.4, 8.7, and 8.8).

**OS\_Bufx** (32 bits, bytes 32-35):

– In *Request\_To\_Receive Operations*: the Buffer Index at the Originating Source or the high-order portion of a 64-bit concatenated address (see 4.4.8);

– In *Data Operations*: the most-significant bytes of the Opaque data (see 4.4.9).

**OS\_Offset** (32 bits, bytes 36-39):

– In *Request\_To\_Receive Operations*: the Originating Source's Offset within a Bufx, or the low-order portion of a 64-bit concatenated address (see 4.4.8);

– In *Clear\_To\_Send Operations*: the Final Destination's initial Offset value (see 4.4.7);

– In *Data Operations*: the least-significant bytes of the Opaque data (see 4.4.9).

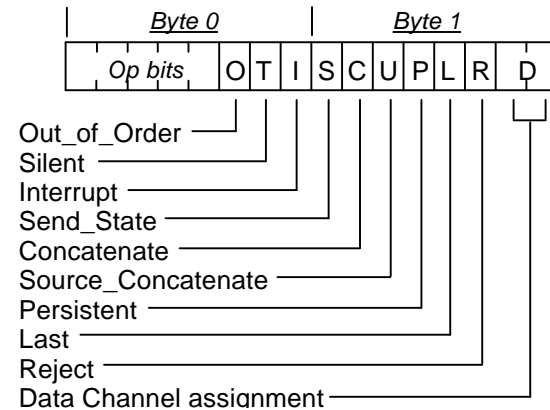
## 6.2 Scheduled Transfer flags

Figure 8 summarizes the flags, and shows their relative position. The flag functions are detailed below for the case where the bit = 1.

**Out\_of\_Order** (b'1xxxxxxxxx') = The end device is able to send and receive Blocks in any order.

**Silent** (b'1xxxxxxxxx') = Requests silent delivery of a Data Operation. For Control Operations the Silent flag shall be ignored (i.e., transmitted as zero, but not checked at the receiver). For Data Operations with Silent = 1 the data transfer to the Destination OS\_Bufx is

carried out normally, but the Schedule Header shall not be delivered to any upper-layer entity. This provides the basis for remote memory write semiotics where the intent is to modify the contents of a remote memory without executing software in the Destination host computer.



**Figure 8 – Flags summary**

**Interrupt** (b'xx1xxxxxxxx') = Requests that a signal or interrupt be generated and delivered to the appropriate upper-layer entity. The Interrupt flag is independent of the Silent flag, i.e., Interrupt = 1 calls for a signal whether or not Silent = 1. (See 4.5.5.)

NOTE 1 – The Silent and Interrupt flags together provide for three delivery modes for Data Operations: silent, polled, or interrupt-driven. If Silent = 1, the data payloads are delivered silently. If Silent = 0, then the upper-layer entity is informed by the same means used for all other Schedule Headers. This mode is suitable for polled interfaces. If Interrupt = 1, then a signal is delivered.

**Send\_State** (b'xxx1xxxxxxxx') = Requests that the Final Destination respond with a **Request\_State\_Response** upon successful receipt of this STU, or Operation, by the higher-layer protocol. **Send\_State is always valid on Control Operations.** For Send\_State to be valid on a Data Operation, either Interrupt = 1 or Silent = 0 must be true.

**Concatenate** (b'xxxx1xxxxx') = Use Bufx as the high-order portion of a 64-bit Final Destination address rather than as a Buffer Index (see 4.3.7).

**Source\_Concatenate** (b'xxxxx1xxxx') = Use OS\_Bufx as the high-order portion of a 64-bit Originating Source address rather than as a Buffer Index (see 4.3.8).

**Persistent** (b'xxxxxx1xxxx') = Retain the Final Destination's buffers (see 4.3.9).

**Last** (b'xxxxxxxx1xxx') = The last STU of a Block.

**Reject** (b'xxxxxxxx1xx') = The request (i.e., Request\_Port, Request\_To\_Send, or Request\_To\_Receive) has been rejected.

**Data Channel assignment:** The Data Channel to be used to carry Data Operations. The Data Channel value is assigned in a Request\_To\_Send Operation and is the Data Channel to be used for Data Operations associated with this Transfer.

b'xxxxxxxx01' = Data Channel 1

b'xxxxxxxx10' = Data Channel 2

b'xxxxxxxx11' = Data Channel 3

The maximum STU size sent on Data Channels 1 and 2 shall be  $2^{17}$  bytes (i.e., 128 Kbytes). The maximum STU size sent on Data Channel 3 shall be  $2^{31}$  bytes (i.e., 2 gigabytes).

NOTE 2 – Data Channel assignment value b'00' is reserved.

## 7 Virtual Connection management

In this clause, a Virtual Connection is setup between two Ports (see 4.3.2), called the A-Port and B-Port. The device that initiates the Virtual Connection is called device A, and the device at the other end is called device B.

In addition to the Port values, each Port shall assign and associate a Key value (A-Key and B-Key) with the Virtual Connection (see 4.3.3). Other parameters exchanged during the Virtual Connection setup include Buffer sizes (A-BuFSIZE and B-BuFSIZE, see 4.3.4), maximum STU sizes (Max-STU, see 4.3.5), and the number of available Slots (A-Slots and B-Slots, see 4.3.6). The end devices also inform each other of their capability to support Concatenate (see 4.3.7), Source\_Concatenate (see 4.3.8), Persistent (see 4.3.9), and out-of-order Block delivery (see 4.4.4).

The Operations used to setup and tear down Virtual Connections are detailed below and summarized in table 2. Only the fields used in each Operation are listed; all of the other Schedule Header fields shall be transmitted as zeros. While a particular field usually carries the parameter of the same name, fields sometimes carry other parameter values. In the Operations below, the specific parameter used in the Operation is listed first, and if it is not carried in the field of the same name, then the field name is included in square brackets.

### 7.1 Request\_Port

Request\_Port shall be used to set up a Virtual Connection between end device A and end device B.

Semantics – Request\_Port (

Op,

Flags,

A-Slots [S\_count],

B-Port [R\_Port],

A-Port [S\_Port],

A-BuFSIZE [Bufx],

A-Key [Offset],

A-Max-STU [T\_len],

EtherType [B\_num] )

Op = x'01'

Flags (see 6.2) shall specify the Out\_of\_Order, Concatenate, Source\_Concatenate, and Persistent flags. A value of 1 shall indicate that A supports that feature. The appropriate value for the Interrupt flag shall also be carried (see 4.5.5).

A-Slots, carried in the S\_count field, shall specify the maximum number of Slots allocated in A for this Virtual Connection (see 4.3.6).

B-Port, carried in the R\_Port field, shall specify B's logical Port value for this Virtual Connection. B-Port may be either the well-known Port (B will assign the Port value), or a peer Port, that provides the service (see 4.3.2).

A-Port, carried in the S\_Port field, shall specify A's logical Port value for this Virtual Connection (see 4.3.2).

A-BuFSIZE, carried in the Bufx field, shall specify A's buffer size (see 4.3.4).



A-Key, carried in the Offset field, shall specify A's Key value for this Virtual Connection (see 4.3.3).

A-Max-STU, carried in the T\_len field, shall be  $\leq$  A-BuFSIZE when sent by A (see 4.3.5). The A-Max-STU value received by B shall be used by B as the maximum size of STUs (Max-STU) to be sent from B to A on this Virtual Connection. (See 4.3.5.)

EtherType, carried in the B\_num field, shall be a value that characterizes the ULP data payloads that will be exchanged on this Virtual Connection (see 4.3.2).

Issued – By device A.

Effect – If it accepts the request, then end device B shall establish a Virtual Connection and shall reply with a Request\_Port\_Response Operation. If rejected, then end device B shall respond with Reject = 1 in the Request\_Port\_Response (see 4.5.3).

## 7.2 Request\_Port\_Response

Request\_Port\_Response shall inform end device A whether the Virtual Connection was accepted or not. If accepted, the parameters associated with this Virtual Connection are passed to A.

Semantics – Request\_Port\_Response (

- Op,
- Flags,
- B-Slots [S\_count],
- A-Port [R\_Port],
- B-Port [S\_Port],
- A-Key [Key],
- B-BuFSIZE [Bufx],
- B-Key [Offset],
- B-Max-STU [T\_len] )

Op = x'02'

Flags (see 6.2) shall specify the Out\_of\_Order, Concatenate, Source\_Concatenate, and Persistent flags. A value of 1 shall indicate that B supports that feature. The appropriate value for the Reject and Interrupt flags shall also be carried (see 4.5.3 and 4.5.5).

B-Slots, carried in the S\_count field, shall specify the maximum number of Slots allocated in B for this Virtual Connection (see 4.3.6).

A-Port, carried in the R\_Port field, shall be the same as the A-Port value in the Request\_Port Operation (see 4.3.2).

B-Port, carried in the S\_Port field, shall specify B's logical Port value for this Virtual Connection (see 4.3.2).

A-Key, carried in the Key field, shall be the Key value assigned by A in the Request\_Port Operation (see 4.3.3).

B-BuFSIZE, carried in the Bufx field, shall specify B's buffer size (see 4.3.4).

B-Key, carried in the Offset field, shall specify B's Key value assigned for this Virtual Connection (see 4.3.3).

B-Max-STU, carried in the T\_len field, shall be  $\leq$  B-BuFSIZE when sent by B (see 4.3.5). The B-Max-STU value received by A shall be used by A as the maximum size of STUs (Max-STU) to be sent from A to B on this Virtual Connection. (See 4.3.5.)

Issued – By B in response to a Request\_Port.

Effect – End device A has been assigned a logical Port on end device B. The Ports, Keys, buffer sizes, maximum STU size, and maximum number of Slots have been exchanged, and a Virtual Connection has been established. Note that the Virtual Connection is bi-directional in that either A or B may initiate a Scheduled Transfer. Multiple Scheduled Transfers may occur over a single Virtual Connection, and the Scheduled Transfers can be either writes or reads. If End Device B will not support the Concatenate or Source\_Concatenate features indicated by the associated bit = 1 in the Request\_Port Operation, then Reject = 1 and the non-supported feature bit = 0, shall be included in the Request\_Port\_Response.

## 7.3 Port\_Teardown

Port\_Teardown shall terminate the Virtual Connection and may be issued by either end device A or end device B. The Port\_Teardown sequence uses a three-way handshake consisting of Port\_Teardown, Port\_Teardown\_ACK, and Port\_Teardown\_Complete.

*Open Issue – A state table describing the 3-way handshake will be included in a Normative annex.*

#### Semantics – Port\_Tearardown (

Op,  
Flags,  
R\_Port,  
S\_Port,  
Key)

Op = x'03'

Flags shall contain the appropriate value for the Interrupt flag (see 4.5.5).

R\_Port shall contain the value associated with the receiver of the Operation, e.g., R\_Port = B-Port when the Port\_Tearardown is issued by A (see 4.3.2).

S\_Port shall contain the value associated with the sender of the Operation, e.g., S\_Port = A-Port when the Port\_Tearardown is issued by A (see 4.3.2).

Key shall contain the Key value associated with the receiver of the Operation, e.g., Key = B-Key when the Port\_Tearardown is issued by A (see 4.3.3).

Issued – By either side, i.e., end device A or end device B, of the Virtual Connection. The sender should only issue a Port\_Tearardown when the Transfers are complete or appear to be stalled.

Effect – The receiver should release any buffers associated with this Virtual Connection, but shall retain the Port and Key values for use in further Port\_Tearardown Operations. The receiver shall also respond with a Port\_Tearardown\_ACK.

### 7.4 Port\_Tearardown\_ACK

Port\_Tearardown\_ACK shall be used to acknowledge receipt of a Port\_Tearardown.

#### Semantics – Port\_Tearardown\_ACK (

Op,  
Flags,  
R\_Port,  
S\_Port,  
Key)

Op = x'04'

Flags shall contain the appropriate value for the

Interrupt flag (see 4.5.5).

R\_Port shall contain the value associated with the receiver of the Operation, e.g., R\_Port = B-Port when the Port\_Tearardown\_ACK is issued by A (see 4.3.2).

S\_Port shall contain the value associated with the sender of the Operation, e.g., S\_Port = A-Port when the Port\_Tearardown\_ACK is issued by A (see 4.3.2).

Key shall contain the Key value associated with the receiver of the Operation, e.g., Key = B-Key when the Port\_Tearardown\_ACK is issued by A (see 4.3.3).

Issued – By the receiver of a Port\_Tearardown Operation after releasing this Virtual Connection's buffers.

Effect – The receiver should release any buffers associated with this Virtual Connection, but shall retain the Port and Key values for use in further Port\_Tearardown Operations. The receiver shall also respond with a Port\_Tearardown\_Complete.

### 7.5 Port\_Tearardown\_Complete

Port\_Tearardown\_Complete shall be used to complete a three-way handshake, acknowledging that the actions associated with a Port\_Tearardown have been completed.

#### Semantics – Port\_Tearardown\_Complete (

Op,  
Flags,  
R\_Port,  
S\_Port,  
Key)

Op = x'05'

Flags shall contain the appropriate value for the Interrupt flag (see 4.5.5).

R\_Port shall contain the value associated with the receiver of the Operation, e.g., R\_Port = B-Port when the Port\_Tearardown\_Complete is issued by A (see 4.3.2).

S\_Port shall contain the value associated with the sender of the Operation, e.g., S\_Port = A-Port when the Port\_Tearardown\_Complete is issued by A (see 4.3.2).

Key shall contain the Key value associated with the receiver of the Operation, e.g., Key = B-Key when the Port\_Teardown\_Complete is issued by A (see 4.3.3).

Issued – By the receiver of a Port\_Teardown\_ACK Operation.

Effect – After the Op\_timeout expires twice, both the sender and receiver shall release the Virtual Connection's Port and Key values. The delay allows for lost or damaged Port\_Teardown Operations to be re-issued.

## 8 Data movement

The Operations used for Scheduled Transfers are detailed below and summarized in table 3. All of the Scheduled Transfer data transfer Operations use the R\_Port, S\_Port, A-Key, and B-Key values that were assigned during the Virtual Connection setup (see 7.1 and 7.2). When end device A issues the Operation:

R\_Port = B-Port  
S\_Port = A-Port  
Key = B-Key

Likewise, when end device B issues the Operation:

R\_Port = A-Port  
S\_Port = B-Port  
Key = A-Key

For clarity and brevity, these values are not discussed in the individual Operations. All other Schedule Header fields that are not listed in a specific Operation shall be transmitted as zeros. While a particular field usually carries the parameter of the same name, fields sometimes carry other parameter values. In the Operations below, the specific parameter used in the Operation is listed first, and if it is not carried in the field of the same name, then the field name is included in square brackets.

### 8.1 Request\_To\_Send

Request\_To\_Send asks that space be allocated, and authorization be given, for a Transfer. Request\_To\_Send is issued by the Originating

Source to specify the number of data bytes to be sent from the Originating Source to the Final Destination. In addition, the Originating Source shall specify whether 64-bit address or Buffer Indexes are used, whether the Final Destination's buffer should be persistent or discarded after a Block, and the Data Channel assignment for the data transfer. Note that the end device on either end of the Virtual Connection may issue a Request\_To\_Send. A Request\_To\_Send, with Persistent = 1, is also used to set up and expose memory for Request\_To\_Receive Operations.

Semantics – Request\_To\_Send (

Op,  
Flags,  
R\_Port,  
S\_Port,  
Key,  
S\_id,  
T\_len [T\_len,B\_num])

Op = x'06'

Flags (see 6.2) shall specify the Concatenate, Persistent, and Data Channel assignment flags. Concatenate or Persistent shall only = 1 if the corresponding flag was set = 1 by the other end during the Virtual Connection setup (see 7.1 and 7.2) and the function is desired for this Operation. The appropriate value for the Interrupt flag shall also be carried (see 4.5.5).

S\_id shall be the Originating Source's Transfer identifier (see 4.4.2) used to identify this Transfer. The Final Destination shall use this value as the R\_id parameter when replying to the Originating Source concerning this Transfer.

T\_len, carried in the concatenation of the T\_len and B\_num fields, shall specify the total number of data payload bytes in the Transfer or that the size of the Transfer is unlimited (see 4.4.3).

Issued – By the Originating Source after a Virtual Connection has been established.

Effect – If accepted, the Final Destination shall set up to receive the data and then reply back to the Originating Source with the associated parameters. If rejected, the Final Destination shall reply with Reject = 1 in a Request\_State\_Response Operation (see 4.5.3).

## 8.2 Request\_To\_Send\_Response

Request\_To\_Send\_Response shall inform the Originating Source whether the Transfer was accepted or not. If accepted, the Request\_To\_Send\_Response specifies the Transfer identifier (see 4.4.2) assigned by this end (i.e., the Final Destination) for this Transfer and the number of STUs per Block (see 4.4.6). A Request\_To\_Send\_Response does not give the Originating Source permission to start sending; that comes from a Clear\_To\_Send. A Clear\_To\_Send may be used instead of a Request\_To\_Send\_Response if the Final Destination is able to immediately accept the data.

Semantics – Request\_To\_Send\_Response (

Op,  
Flags,  
Blocksize [S\_count],  
R\_Port,  
S\_Port,  
Key,  
R\_id,  
S\_id)

Op = x'07'

Flags (see 6.2) shall specify the Concatenate flag. Concatenate shall only = 1 if Concatenate was set = 1 by the remote end device during the Virtual Connection setup (see 7.1 and 7.2). The appropriate value for the Reject and Interrupt flags shall also be carried (see 4.5.3 and 4.5.5).

Blocksize, carried in the S\_count field, shall specify the **Block size** (see 4.4.5).

**NOTE – Blocksize is the number of bytes in a Block, expressed as a power of 2 (see 4.4.5). Blocksize is not the number of STUs in a Block.**

R\_id shall be the Transfer identifier (see 4.4.2) assigned by the Originating Source in the Request\_To\_Send Operation.

S\_id shall be the Transfer identifier (see 4.4.2) used by the Final Destination to identify this Transfer. The Originating Source shall use this value as the R\_id parameter when replying to the Final Destination concerning this Transfer.

Issued – By the Final Destination.

Effect – The Originating Source shall segment the Transfer into Blocks and STUs for transmission.

## 8.3 Request\_To\_Receive

Request\_To\_Receive, issued by the Final Destination, asks for a single Block of data to be sent from a previously allocated location in the Originating Source. The Request\_To\_Receive specifies the number of data bytes to be sent from the Originating Source to the Final Destination and whether 64-bit addresses or Buffer Indexes are used at the Originating Source, at the Final Destination, or at both. A Request\_To\_Receive transfers a single Block; there is no notion of a multi-Block Request\_To\_Receive data movement.

When a Request\_To\_Receive is issued, it is assumed that the ULPs on both end devices had previously allocated resources for the entire Transfer through a previous Request\_To\_Send Operation. Note that the device at either end of the Virtual Connection may issue a Request\_To\_Receive.

**Request\_To\_Receive may be used in conjunction with Persistent memory. The Request\_To\_Receive initiator must request a remote Persistent memory region by issuing a Request\_To\_Send Operation with Persistent = 1. If accepted, an (optional) Request\_To\_Send\_Response will be returned with Persistent = 1. This will be followed by a Clear\_To\_Send Operation which establishes the memory region dedicated for the Transfer. The RTS/CTS handshake establishes an R\_ID, S\_ID pair which together identify the "Transfer". The Persistent memory for the Transfer remains available until the Transfer is terminated by either an END/END\_ACK exchange or Port\_Teardown sequence.**

**While a Persistent Transfer is active, it is available for an unlimited number of Data Operations or Request\_To\_Receive Operations as long as the number of outstanding Operations at any time falls within the limits established by the Slot mechanism (section 4.3.6).**

## Semantics – Request\_To\_Receive (

Op,  
 Flags,  
 R\_Port,  
 S\_Port,  
 Key,  
 R\_id,  
 S\_id,  
 Bufx,  
 Offset,  
 T\_len [T\_len,B\_num],  
 OS\_Bufx,  
 OS\_Offset )

Op = x'08'

Flags (see 6.2) shall specify the Concatenate and Source\_Concatenate flags. Source\_Concatenate shall only = 1 if the corresponding flag was set by the remote end device during the Virtual Connection setup (see 7.1 and 7.2) and the function is desired for this Operation. The appropriate value for the Interrupt flag shall also be carried (see 4.5.5).

R\_id shall be the Transfer identifier (see 4.4.2) assigned by the Originating Source in the Request\_To\_Send Operation.

S\_id shall be the Transfer identifier (see 4.4.2) used by the Final Destination to identify this Transfer. The Originating Source shall use this value as the R\_id parameter when replying to the Final Destination concerning this Transfer.

Bufx shall specify the initial Buffer Index in the Final Destination where the data will be placed (see 4.4.7).

Offset is a value that the Final Destination must receive with the first STU of the Block so that the data can be properly placed in the Final Destination's memory (see 4.4.7).

T\_len, carried in the concatenation of the T\_len and B\_num fields, shall specify the total number of data payload bytes in the Transfer (see 4.4.3).

OS\_Bufx shall specify the Originating Source's Buffer Index (see 4.4.8).

OS\_Offset shall specify the Originating Source's offset value (see 4.4.8).

Issued – By the Final Destination.

Effect – If accepted, the Originating Source shall

send the specified Block of data. If rejected, the Originating Source shall reply with Reject = 1 in a Request\_To\_Receive\_Response (see 4.5.3).

#### 8.4 Request\_To\_Receive\_Response

Request\_To\_Receive\_Response, issued by the Originating Source, tells the Final Destination that the Request\_To\_Receive that it issued has been rejected.

#### Semantics – Request\_To\_Receive\_Response (

Op,  
 Flags,  
 R\_Port,  
 S\_Port,  
 Key,  
 R\_id,  
 S\_id)

Op = x'09'

Flags (see 6.2) shall specify the Reject and Interrupt flags.

R\_id shall be the Transfer identifier (see 4.4.2) assigned by the Originating Source in the Request\_To\_Send Operation.

S\_id shall be the Transfer identifier (see 4.4.2) used by the Final Destination to identify this Transfer.

Issued – By the Originating Source.

Effect – The Request\_To\_Receive has been rejected (see 4.5.3).

#### 8.5 Clear\_To\_Send

Clear\_To\_Send shall be used to give the Originating Source permission to send one Block. Clear\_To\_Send may also be used to request retransmission of a Block from systems that are capable of retransmission.

#### Semantics – Clear\_To\_Send (

Op,  
Flags,  
Blocksize [S\_count],  
R\_Port,  
S\_Port,  
Key,  
R\_id,  
S\_id,  
Bufx,  
Offset,  
B\_num,  
I\_Offset [OS\_Offset])

Op = x'0A'

Flags (see 6.2) shall specify the Concatenate and Source\_Concatenate flags. These flags shall be the same value as in the Request\_To\_Send Operation that initiated this Transfer. The appropriate value for the Interrupt flag shall also be carried (see 4.5.5).

*Open Issue – Usage of the Concatenate and Source\_Concatenate flags in the data movement operations needs to be clarified.*

Blocksize, carried in the S\_count field, shall specify the **Block size** (see 4.4.5).

R\_id shall be the Transfer identifier (see 4.4.2) assigned by the remote end (the Originating Source) of the Virtual Connection.

S\_id shall be the Transfer identifier (see 4.4.2) assigned by this end (the Final Destination) of the Virtual Connection.

Bufx shall specify the initial Buffer Index in the Final Destination where the data will be placed (see 4.4.7).

Offset is a value that the Final Destination must receive with the first STU of a Block so that the data can be properly placed in the Final Destination's memory (see 4.4.7).

B\_num shall be the Block number being given permission to be transmitted (see 4.4.4).

I\_Offset, carried in the OS\_Offset field, shall specify the Offset value associated with the first Block of the Transfer (see 4.4.7).

Issued – By the Final Destination.

Effect – The Originating Source shall send the specified Block.

## 8.6 Data

A Data Operation sends an STU of a Block from the Originating Source to the Final Destination. No STU shall be larger than the maximum STU size determined during the Virtual Connection setup (see 7.2).

#### Semantics – Data (

Op,  
Flags,  
S\_count,  
R\_Port,  
S\_Port,  
Key,  
R\_id,  
S\_id,  
Bufx,  
Offset,  
Sync [T\_len],  
B\_num,  
Opaque [OS\_Bufx],  
Opaque [OS\_Offset])

Op = x'0B'

Flags (see 6.2) shall specify the Interrupt, Silent, Send\_State, Concatenate, Last, and Data Channel assignment flags (see 6.2). Concatenate shall be the same value as in the Request\_To\_Send or Request\_To\_Receive that initiated this Transfer. Send\_State may be sent with any STU of a Block. The **Request\_State\_Response** Control Operation associated with this request shall be issued after processing this STU when Send\_State = 1. The sender shall copy (in this Data Operation) the Data Channel assignment flags supplied in the corresponding Request\_To\_Send Operation; the value is a do not care at the receiver.

S\_count shall be the STU number (see 4.4.6).

R\_id shall be the Transfer identifier (see 4.4.2) assigned by the other end (the Final Destination) of the Virtual Connection.

S\_id shall be the Transfer identifier (see 4.4.2) assigned by this end of the Virtual Connection. Note that if this is the first STU associated with a Request\_To\_Receive Operation, then this Transfer identifier (see 4.4.2) is being assigned by the Originating Source and shall be used by the Final Destination as the R\_id parameter when replying to the Originating Source



concerning this Transfer.

Bufx shall be the Buffer Index at the Final Destination (see 4.4.7).

Offset shall be the Final Destination's offset within a Bufx (see 4.4.7).

Sync, carried in the T\_len field, shall be a value assigned by the Originating Source to synchronize the current view of the number of empty Slots in the Final Destination (see 4.3.6).

B\_num shall be the number of the Block that this STU is a part of.

Opaque data, carried in the OS\_Bufx and OS\_Offset fields, shall be as specified in 4.4.9.

Issued – By the Originating Source.

Effect – The Final Destination shall place the STU data in the memory area pointed to by Bufx and offset by the Offset value. The Final Destination shall only accept data into pre-allocated buffer regions. The Final Destination is responsible for ensuring that all of the Blocks of a Transfer are received. The actions to be taken if a Block is missing are beyond the scope of this standard.

## 8.7 Request\_State

Request\_State is used to request that the remote end device provide its current number of empty Slots for Schedule Headers, the Block number associated with the last set of contiguously good data received, and whether the named Block was received correctly.

Semantics – Request\_State (

Op,  
Flags,  
R\_Port,  
S\_Port,  
Key,  
R\_id,  
S\_id,  
Sync [T\_len],  
B\_num )

Op = x'0C'

Flags shall contain the appropriate value for the Interrupt flag (see 4.5.5).

R\_id shall be the Transfer identifier (see 4.4.2) assigned by the remote end device of this Virtual Connection. R\_id = x'0000' means that the receiver shall not look for a current Transfer and only return the current number of empty Slots for this Virtual Connection.

S\_id shall be the Transfer identifier (see 4.4.2) assigned by this end of the Virtual Connection. If R\_id = x'0000', then S\_id shall also be x'0000'.

Sync, carried in the T\_len field, shall be a value assigned by the local end device (sender) to synchronize the current view of the number of empty Slots in the remote end device (receiver). (See 4.3.6.)

B\_num shall indicate the Block number being queried. B\_num = x'FFFFFFFF' indicates that the sender does not care about the status of any particular Block.

Issued – By an end device that needs state information from the remote end device of the Virtual Connection. The sender may not have received the Request\_State\_Response that it expected from a Data Operation and can send a Request\_State to recover from a lost or damaged Request\_State\_Response.

Effect – The receiver shall reply with a Request\_State\_Response.

## 8.8 Request\_State\_Response

Request\_State\_Response shall be used to indicate the number of empty Slots in this Port of the Virtual Connection (see 4.3.6). Request\_State\_Response may also indicate the highest numbered contiguous Block received correctly and whether the Block indicated in the B\_num parameter was received correctly (see 4.5.2).

Semantics – Request\_State\_Response (

Op,  
Flags,  
C-Slots [S\_count],  
R\_Port,  
S\_Port,  
Key,  
R\_id,  
S\_id,  
B\_seq [Offset],  
Sync [T\_len],  
B\_num)

Op = 'x'0D'

Flags (see 6.2) shall specify Reject = 1 if the immediately previous Request\_To\_Receive was rejected (see 4.5.3). The appropriate value for the Interrupt flag shall also be carried (see 4.5.5).

C-Slots, carried in the S\_count field, shall indicate the sender's view of the number of empty Slots it has available for additional Operations on this Virtual Connection. (See 8.6.) C-Slots = 'x'FFFF' (i.e., -1) shall indicate that this end device does not implement the Slots mechanism for Operations flow control.

R\_id shall echo the S\_id value in the Request\_State or Data Operation that triggered this Request\_State\_Response.

S\_id shall echo the R\_id value in the Request\_State or Data Operation that triggered this Request\_State\_Response. S\_id = 'x'0000' shall mean that the B\_seq and B\_num parameters are meaningless.

B\_seq, carried in the Offset field, shall indicate the highest numbered contiguous Block received correctly. B\_seq = 'x'FFFFFFFF' shall indicate that no Transfers are in progress or no Blocks have been received.

Sync, carried in the T\_len field, is echoed from the Request\_State, or Data Operation with Send\_State = 1, that initiated this Request\_State\_Response Operation (see 4.3.6).

B\_num shall echo the Block number, carried in the B\_num field of the Data Operation or the Request\_State Operation, if the indicated Block was received correctly. If the indicated Block has not been correctly received, then B\_num shall contain 'x'FFFFFFFF'. The Sync value can be used by the receiving end to identify the

Operation containing the B\_num being queried.  
Issued – It is intended that a Request\_State\_Response be issued by an end device's ULP after receiving Send\_State = 1 in a Data Operation, or after receiving a Request\_State Operation, or to reject an Operation.

Effect – State information is passed to the other end of the Virtual Connection.

## 8.9 END

END allows either end of the Virtual Connection to terminate a Scheduled Transfer before it has completed and to terminate a Scheduled Transfer of unlimited size.

Semantics – END (

Op,  
Flags,  
R\_Port,  
S\_Port,  
Key,  
R\_id  
S\_id)

Op = 'x'0E'

Flags shall contain the appropriate value for the Interrupt flag (see 4.5.5).

R\_id shall be the Transfer identifier (see 4.4.2) assigned by the other end of the Virtual Connection.

S\_id shall be the Transfer identifier (see 4.4.2) assigned by this end of the Virtual Connection. This S\_id value shall not be reused until an END\_ACK is received.

Issued – By the Originating Source or the Final Destination.

Effect – A Final Destination receiving an END shall stop sending Control Operations associated with this Scheduled Transfer. An Originating Source receiving an END shall stop sending Control Operations and STUs associated with this Scheduled Transfer. An END kills a Scheduled Transfer, but shall not affect the Virtual Connection carrying the Scheduled Transfer.



**Table 2 – Virtual Connection Operations summary between end devices A and B**

	Op	Flags	S_count	R_Port	S_Port	Key	Bufx	Offset	T_len	B_num
RQP	x'01'	<b><i>OIUCP</i></b>	<b><i>A-Slots</i></b>	<b><i>B-Port</i></b>	<b><i>A-Port</i></b>	*	<b><i>A-Bufsize</i></b>	<b><i>A-Key</i></b>	<b><i>A-Max-STU</i></b>	<b><i>EtherType</i></b>
RQPR	x'02'	<b><i>OIUCPR</i></b>	<b><i>B-Slots</i></b>	A-Port	<b><i>B-Port</i></b>	A-Key	<b><i>B-Bufsize</i></b>	<b><i>B-Key</i></b>	<b><i>B-Max-STU</i></b>	*
PT	x'03'	<b><i>I</i></b>	*	R_Port	S_Port	R-Key	*	*		*
PTA	x'04'	<b><i>I</i></b>	*	R_Port	S_Port	R-Key	*	*		*
PTC	x'05'	<b><i>I</i></b>	*	R_Port	S_Port	R-Key	*	*		*
NOTES – 1 – Operation abbreviations: PT = Port_Teardown PTA = Port_Teardown_ACK PTC = Port_Teardown_Complete RQP = Request_Port RQPR = Request_Port_Response 2 – Flag abbreviations are: O = Out_of_Order, I = Interrupt, U = Source_Concatenate, C = Concatenate, P = Persistent, R = Reject 3 – R-Key = Key value the receiver binds to, e.g., R-Key = A-Key when Operation issued by device B. 4 – R_Port = Port number in device receiving the Operation, e.g., R_Port = A-Port when issued by device B. 5 – S_Port = Port number in device sending the Operation, e.g., S_Port = B-Port when issued by device B. 6 – The Schedule Header fields that are not shown shall be transmitted as zeros. SYMBOLS - * = Unused value, transmit as 0 Values in bold italics are assigned by the specific Operation and may be used by later Operations										

**8.10 END\_ACK**

END\_ACK confirms that the sending end device has seen, and acted on, the END.

Semantics – END\_ACK (

Op,  
Flags,  
R\_Port,  
S\_Port,  
Key,  
R\_id,  
S\_id)

Op = **x'0F'**

Flags shall contain the appropriate value for the Interrupt flag (see 4.5.5).

R\_id shall be the Transfer identifier (see 4.4.2) assigned by the remote end device of this Virtual Connection.

S\_id shall be the Transfer identifier (see 4.4.2) assigned by this end of the Virtual Connection. This S\_id value should not be immediately reused to avoid aliasing.

Issued – By the end of the Virtual Connection that received the END Operation.

Effect – Acknowledgment that the Scheduled Transfer has been terminated.

**Table 3 – Data transfer and status Operations summary between end devices S and R**

	Op	Flags	S_count	R_id	S_id	Bufx	Offset	T_len	B_num	OS_Bufx	OS_Offset
RTS	x'06'	<b>ICPD</b>	*	*	<b>S_id</b>	*	*	<b>T_len</b>		*	*
RTSR	x'07'	<b>ICR</b>	<b>Blocksize</b>	R_id	<b>S_id</b>	*	*	*	*	*	*
RTR	x'08'	<b>ICU</b>	*	R_id	S_id	<b>Bufx</b>	<b>Offset</b>	<b>T_len</b>		<b>OS_Bufx</b>	<b>OS_Offset</b>
<b>RTRR</b>	<b>x'09'</b>	<b>IR</b>	*	<b>R_id</b>	<b>S_id</b>	*	*	*	*	*	*
CTS	x'0A'	<b>ICU</b>	<b>Blocksize</b>	R_id	S_id	<b>Bufx</b>	<b>Offset</b>	*	<b>B_num</b>	*	<b>I_Offset</b>
Data	x'0B'	<b>ITSCLD</b>	<b>S_count</b>	R_id	S_id	<b>Bufx</b>	<b>Offset</b>	<b>Sync</b>	B_num	<b>Opaque</b>	<b>Opaque</b>
RS	x'0C'	<b>I</b>	*	R_id	S_id	*	*	<b>Sync</b>	<b>B_num</b>	*	*
<b>RSR</b>	<b>x'0D'</b>	<b>IR</b>	<b>C-Slots</b>	R_id	S_id	*	<b>B_seq</b>	Sync	B_num	*	*
END	x'0E'	<b>I</b>	*	R_id	S_id	*	*	*	*	*	*
ENDA	x'0F'	<b>I</b>	*	R_id	S_id	*	*	*	*	*	*

**NOTES –**

## 1 – Operation abbreviations:

CTS = Clear\_To\_Send

ENDA = END\_ACK

RS = Request\_State

**RSR = Request\_State\_Response**

RTR = Request\_To\_Receive

**RTRR = Request\_To\_Receive\_Response**

RTS = Request\_To\_Send

RTSR = Request\_To\_Send\_Response

## 2 – Flag abbreviations : I = Interrupt, T = Silent, S = Send\_State, C = Concatenate, U = Source\_Concatenate, P = Persistent, L = Last STU of Block, R = Reject, D = Data Channel assignment

## 3 – R\_id = Transfer identifier in device receiving the Operation, e.g., R\_id = G\_id when issued by device H.

## 4 – S\_id = Transfer identifier in device sending the Operation, e.g., S\_id = H\_id when issued by device H.

## 5 – Schedule Header parameters that shall be transmitted as assigned in RQP and RQPR Operations:

R\_Port = Port number of the device receiving the Operation

S\_Port = Port number of the device sending the Operation

Key = Key value assigned by the device receiving the Operation

**SYMBOLS -**

\* = Unused value, transmit as 0

Values in bold italics are assigned by the specific Operation and may be used by later Operations

## 9 Error processing

Table 5 is a summary of the logged errors. The logging is on a per-Port basis, and shall be available to the ULP that **is using** the Port. The nature and size of the logs are system dependent.

### 9.1 Operation timeout

Errors other than syntactic errors are manifested as missing Operations, occurring when the underlying physical media discard or damage a transmission (see 4.5.4). Such errors are detected by Op\_timeout, which is system and/or Port dependent. Op\_timeout\_Occurances shall be logged. Example means for determining the Op\_timeout value for a Virtual Connection include:

- a time longer than the measured round-trip time through the software path (use a Request\_State / **Request\_State\_Response** pair to measure on a per-Port basis); or
- a long fixed time period.

Another system and/or Port dependent parameter, Max\_Retry, specifies the maximum number of times to retry an Operation. If Max\_Retry is reached without success, then the Operation is considered to be aborted and control shall be passed to the ULP. Max\_Retry\_Occurances shall be logged.

### 9.2 Operation Pairs

Table 4 lists the Operation pairs – command and response, or response and completion – that shall be retried if the associated response is not received within an Op\_timeout.

In addition to the entries in table 4, **Request\_State\_Response** is a corresponding pair for Data Operations which have Send\_State = 1. If the **Request\_State\_Response** is not received, then the Originating Source may send a Request\_State to obtain the state information.

The ULP in the Final Destination that issues a Clear\_To\_Send, or Request\_To\_Receive, is responsible for timing out these Operations. The ULP may or may not use Op\_timeout to indicate failure.

**Table 4 – Operation pairs guarded by Op\_timeout**

Operation	Response
Request_Port	Request_Port_Response
Port_Teardown	Port_Teardown_ACK
Port_Teardown_ACK	Port_Teardown_Complete
Request_To_Send	Request_To_Send_Response, or Clear_To_Send
Request_To_Receive	Data, or <b>Request_To_Receive_Response</b>
Request_State	<b>Request_State_Response</b>
END	END_ACK

### 9.3 Syntax errors

#### 9.3.1 Undefined Opcode

An undefined Opcode value may occur due to bit errors or if the sending device is using a future superset of the Scheduled Transfer Operations. The Operation shall be discarded, an Undefined\_Opcode\_Error shall be logged, and the Opcode logged in Undefined\_Opcode\_Value.

#### 9.3.2 Unexpected Opcode

Most of the Operations require previous Operations to set up state on each device. If a device receives an out of sequence Opcode (e.g., receiving a Request\_Port\_Response without sending the initiating Request\_Port), the Operation shall be discarded, an Unexpected\_Opcode\_Error shall be logged, and the Opcode logged in Unexpected\_Opcode\_Value.

### 9.4 Virtual Connection errors

#### 9.4.1 Invalid Key or Port

All Operations, excluding Request\_Port, should have a Key (**see 4.3.3**) value that validates the Operation for the Virtual Connection. Operations with an invalid Key shall not be executed, and an Invalid\_Key\_Error shall be logged.

All Operations should have a valid Destination Port value (see 4.3.2). Operations with an invalid Destination Port value shall not be executed, and an Invalid\_Port\_Error shall be logged.

NOTE – Multiple contiguous invalid Key and/or Port values may indicate a problem with the link or a malicious host on the network. The supervising process should be informed.

#### 9.4.2 Slots exceeded

Operations that exceed the number of Slots (see 4.3.6) for the Virtual Connection may not be executed, and a Slots\_Exceeded\_Error shall be logged.

#### 9.4.3 Unknown EtherType

If a Request\_Port Operation contains an unknown EtherType (see 4.3.2), the receiver shall issue a Request\_Port\_Response with the Reject bit set and log an Unknown\_EtherType\_Error.

#### 9.4.4 Illegal Buftime

If a Request\_Port contains a Buftime (see 4.3.4) value that is  $< 8$  or  $> 64$ , (i.e., Buffer size  $< 2^8$  bytes, or  $> 2^{64}$  bytes), then the receiver shall respond with a Request\_Port\_Response with Reject = 1. If a Request\_Port\_Response contains a Buftime value that is  $< 8$  or  $> 64$ , then the receiver shall respond with a Port\_Teardown. In either case, an Illegal\_Buftime\_Error shall be logged.

#### 9.4.5 Illegal STU size

The maximum STU sizes (A-Max-STU and B-Max-STU) were determined during the Virtual Connection setup (see 4.3.5, 7.1 and 7.2). If the received STU is greater than the maximum STU size, then the STU shall be discarded and an Illegal\_STU\_Size\_Error shall be logged.

### 9.5 Scheduled Transfer errors

#### 9.5.1 Invalid S\_id

All Scheduled Transfer Operations, except Request\_To\_Send, should have a valid

Destination id (R\_id) (see 4.4.2) for quickly accessing state information for this Scheduled Transfer. After checking the R\_id, the S\_id should match the stored value for this Transfer. An invalid S\_id shall result in not executing the Operation and logging an Invalid\_S\_id\_Error.

#### 9.5.2 Bad Data Channel specification

During a Request\_To\_Send Operation, the sending device declares the lower layer Data Channel that will carry Data Operations for this Scheduled Transfer. Some Data Channels may not be available for Scheduled Transfers depending on the lower layer (e.g., b'00' is not a valid choice on HIPPI-6400 as it indicates VC0 which is reserved for Control Operations). The receiver shall issue a Request\_To\_Send\_Response with the Reject bit set.

#### 9.5.3 Concatenate not available

If the Virtual Connection did not specify the capability for Concatenate (see 4.3.7) during the Virtual Connection establishment (see 7.1 and 7.2), any Scheduled Transfer Operations on this Virtual Connection with the Concatenate bit set shall not be executed. The Operation shall be rejected (see 4.5.3) and a Concatenate\_Error logged.

#### 9.5.4 Source\_Concatenate not available

If the Virtual Connection did not specify the capability for Source\_Concatenate (see 4.3.8) during the Virtual Connection establishment (see 7.1 and 7.2), any Scheduled Transfer Operations on this Virtual Connection with the Source\_Concatenate bit set shall not be executed. The Operation shall be rejected and a Source\_Concatenate\_Error logged.

#### 9.5.5 Persistent not available

If the Virtual Connection did not specify the capability for Persistent (see 4.3.9) during the Virtual Connection establishment (see 7.1 and 7.2), any Scheduled Transfer Operations on this Virtual Connection with the Persistent bit set shall not be executed. The Operation shall be rejected and a Persistent\_Error logged.

### 9.5.6 Out of Range B\_num, Bufx, Offset, S\_count, or Sync

During the Clear\_To\_Send, Data, and Request\_State\_Response Operations, a Block number (see 4.4.4) may appear that is outside the calculated number of Blocks for the Transfer. If an out of range Block number is encountered, the receiver shall not execute the Operation and shall log an Out\_Of\_Range\_B\_num\_Error.

If a Data Operation contains a Bufx and/or Offset (see 4.4.7) that exceeds the buffer range allocated by the Final Destination for outstanding Clear\_To\_Sends, then the receiver shall not execute the Operation and shall log an Out\_Of\_Range\_Bufx\_Error.

If a Data Operation contains an Offset (see 4.4.7) larger than the buffer size, the receiver shall not execute the Operation and shall log an Oversized\_Offset\_Error.

If a Data Operation contains an S\_count (see 4.4.6) that is not one greater than the previous STU (for this Block), then the STU is out of order. The receiver shall discard the STU and log an Out\_Of\_Order\_STU\_Error.

### 9.5.7 Illegal Blocksize

If a Request\_To\_Send\_Response, or Clear\_To\_Send, contains a Blocksize (see 4.4.5) value that is  $< 8$  or  $> 64$ , (i.e., Block size  $< 2^8$  bytes, or  $> 2^{64}$  bytes), then the receiver should discard the offending Operation and log an Illegal\_Blocksize\_Error.

*Open Issue – This is a new check and needs to be reviewed. Note the use of "should" rather than a "shall".*

### 9.5.8 Request\_To\_Receive problem

The Request\_To\_Receive Operation (see 8.3) must be set up by previous Request\_To\_Send and Clear\_To\_Send Operations. If these Operations have not been received, then the Request\_To\_Receive Operation shall be discarded, a Request\_State\_Response with Reject = 1 sent to the remote end device in response, and a Request\_To\_Receive\_Error logged.

### 9.5.9 Undefined Flag

If a received Operation contains a flag =1 and use of that flag is not defined for that Operation, then the flag shall be ignored and an Improper\_Flag\_Use\_Error logged.

**Table 5 – Summary of logged errors**

Name	Occurs in Operation
Concatenate_Error	RTS, RTSR, RTR, CTS
Illegal_Blocksize_Error	RTSR, CTS
Illegal_Bufsize_Error	RQP, RQPR
Illegal_STU_Size_Error	Data
Improper_Flag_Use_Error	all
Invalid_Key_Error	all except RQP
Invalid_Port_Error	all
Invalid_S_id_Error	all with an R_id
Max_Retry_Occurance	END, PT, PTA, RQP, RS, RTR, RTS
Op_timeout_Occurance	END, PT, PTA, RQP, RS, RTR, RTS
Out_Of_Order_STU_Error	Data
Out_Of_Range_B_num_Error	CTS, Data, RS, SR
Out_Of_Range_Bufx_Error	Data
Oversized_Offset_Error	Data
Persistent_Error	RTS
Request_To_Receive_Error	RTR
Slots_Exceeded_Error	all with Opcode $\geq 6$
Source_Concatenate_Error	RTR, CTS
Undefined_Opcode_Error	not applicable
Undefined_Opcode_Value	not applicable
Unexpected_Opcode_Error	all except RQP
Unexpected_Opcode_Value	all except RQP
Unknown_EtherType_Error	RQP
Operation abbreviations: CTS = Clear_To_Send PT = Port_Teardown PTA = Port_Teardown_ACK RQP = Request_Port RQPR = Request_Port_Response RS = Request_State RTR = Request_To_Receive RTS = Request_To_Send SR = Request_State_Response	

## **Annex A** (normative)

### **Using lower layer protocols**

#### **A.1 HIPPI-6400-PH as the lower layer**

ANSI X3.xxx defines HIPPI-6400-PH, portions of which are repeated here as an aid to the reader. As shown in figure A.1, HIPPI-ST Operations shall be carried over HIPPI-6400-PH with the first eight bytes of the HIPPI-ST Schedule Header occupying the last eight bytes of the HIPPI-6400-PH Header micropacket.

All HIPPI-ST Control Operations shall be carried on HIPPI-6400-PH Virtual Channel VC0. Data Operations shall use Virtual Channel 1, 2, or 3 as specified in the HIPPI-ST Data Channel Assignment flag bits (see 6.2) and carried in a Request\_To\_Send Operation (see 8.1).

HIPPI-ST shall also specify the EtherType value that is placed in the HIPPI-6400-PH MAC Header (see the reference for RFC 1700 in 4.3.2).

M\_len (in the HIPPI-6400-PH MAC Header), specifies the number of bytes following M\_len, exclusive of any padding in the last micropacket. Hence, M\_len will have the following values:

- M\_len = 48 for Control Operations without an optional payload (i.e., 48 = 8 byte IEEE 802.2 LLC/SNAP Header + 40-byte HIPPI-ST Schedule Header);
- M\_len = 80 for Control Operations with optional payload;
- M\_len = (48 + number of user data payload bytes) for Data Operations.

#### **A.2 HIPPI-FP as the lower layer**

ANSI X3.210 defines HIPPI-FP, portions of which are repeated here as an aid to the reader. As shown in figure A.2, HIPPI-ST Operations shall be carried over HIPPI-FP in the D2\_Area. The HIPPI-FP D1\_Area shall not be used. The HIPPI-FP D2\_Offset shall be set to zero. Short bursts shall only be used at the end of a packet, i.e., short first burst is disallowed. Note that D2\_Size = M\_len + 16.

The HIPPI-6400-PH MAC and LLC/SNAP Headers are defined in ANSI X3.xxx, portions of which are repeated here as an aid to the reader. The MAC and LLC/SNAP headers are included to facilitate translation to other protocols. The 48-bit ULA addresses allow address assignment and usage common to other networking technologies.

*Open Issue – Address mapping between HIPPI-800 and HIPPI-6400 devices has not been worked out, in particular for non-channel attached devices.*

*Open Issue – The Introduction says that this document specifies mappings to IPv4, IPv6, and MPI upper-layer protocols. Where is this information going to be placed? Who is going to provide it?*

HIPPI-6400-PH MAC and LLC/SNAP Headers	D_ULA				32-byte HIPPI-6400-PH Type = Header micropacket
	D_ULA (lsb)		S_ULA		
	S_ULA (lsb)				
	M_len				
	DSAP	SSAP	Ctl	Org	
	Org	Org	EtherType		
HIPPI-ST Header (defined in 6.1 and shown here as an aid to the reader)	Op	Flags	S_count		First 32-byte HIPPI-6400-PH Type = Data micropacket
	R_Port		S_Port		
	Key				
	R_id		S_id		
	Bufx				
	Offset				
	T_len				
	B_num				
	OS_Bufx				
	OS-Offset				
HIPPI-ST payload	Optional 32-byte payload (in Control Operations)  or  Up to 2 <sup>31</sup> bytes (2 gigabytes) of HIPPI-ST data payload (i.e., STU) (in Data Operations)				Additional 32-byte HIPPI-6400-PH Type = Data micropacket(s)

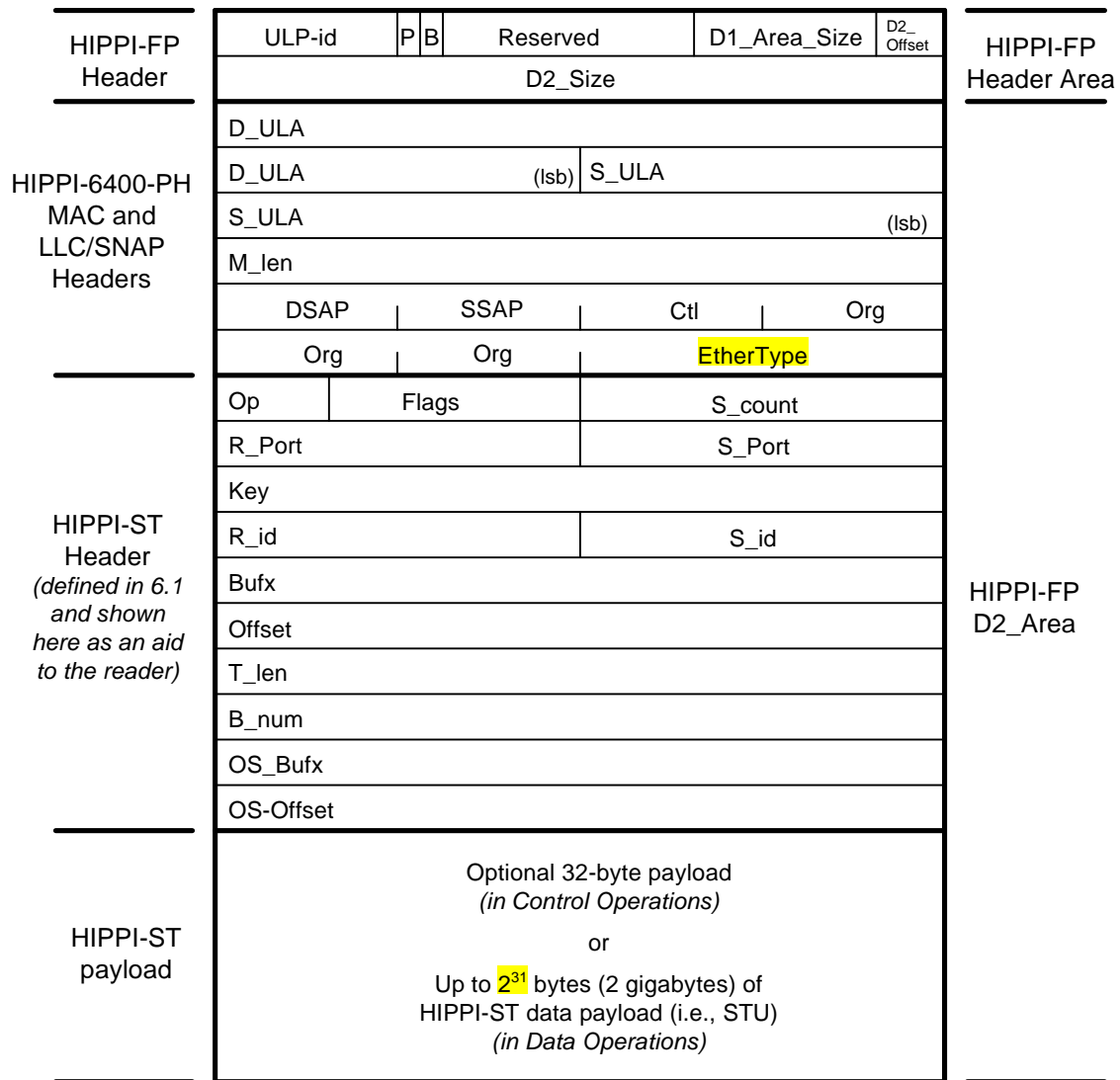
32-byte  
HIPPI-6400-PH  
Type = Header  
micropacket

First 32-byte  
HIPPI-6400-PH  
Type = Data  
micropacket

Additional 32-byte  
HIPPI-6400-PH  
Type = Data  
micropacket(s)

NOTE – Shown as 32-bit worrds

**Figure A.1 – HIPPI-ST Operations carried in HIPPI-6400-PH Messages**



NOTE – Shown as 32-bit words

**Figure A.2 – HIPPI-ST Operations carried in HIPPI-FP packets**



## Annex B (informative)

### HIPPI-ST striping

#### B.1 Striping principles

HIPPI-ST is capable of supporting multiple physical interfaces for a single Transfer (see figures B.1–B.3). This striping capability may be of benefit when a single interface is not able to support required data rates. It may be especially useful where data is moved from many slower interfaces to a single faster interface or vice-versa. It may also be used with multiple interfaces at both the Originating Source and Final Destination.

**The Block is the basic striping unit.** Each Block contains sufficient information to completely identify an individual Transfer and the Block's location within the Transfer. The only difference between striped and non-striped operation is the selection of port MAC addresses to allow concurrent data movement. Striping is performed on a Block, not a STU basis. Otherwise, STU in-order delivery is unlikely to be guaranteed.

There are a few conventions that should be followed to facilitate striping:

- Block sizes (when striping is desirable) must be small enough to support concurrency and allow each channel to have at least one Block to send.
- Sufficient CTS operations should be kept outstanding by a data receiver to allow concurrent Data operations.
- The interface adapter(s) must be capable of handling multiple Blocks simultaneously. This may require communication between interfaces (or their software drivers) within a system.
- The return ULA for each operation is specified by the Source ULA for that operation. HIPPI-ST implementers should not assume that the Source ULA for a given port will remain constant.

#### B.2 Many-to-one striping

When a number of lower bandwidth interfaces are aggregated to communicate with one faster interface (using a translator or bridge), striping the lower bandwidth interfaces together can allow legacy systems to communicate quickly over newer network infrastructures. In this case, action to implement striping is required only on the side of the lower bandwidth interface.

After port set-up, a Transfer is initiated with a RTS operation. An RTS\_Response will be received, either as a discrete message or as part of a CTS. As CTS operations are received, the system with multiple lower bandwidth ports can move a Block of data for each CTS received. As many Blocks can be in transit concurrently as there are ports to send them on and CTS operations authorizing them.

The system receiving these Blocks processes them normally, placing them into memory as their Bufx and Offset values dictate.

#### B.3 One-to-many striping

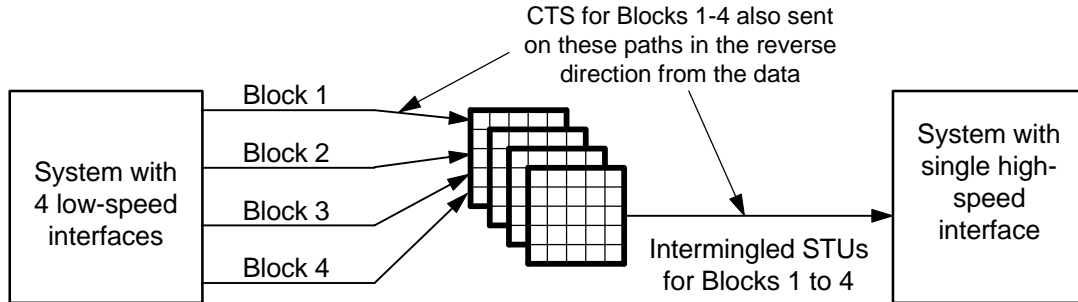
Transfers made from one high-speed interface can also be spread across more than one low-speed interface without any special action on the part of the high speed system.

After port setup, the Transfer is initiated with a RTS operation from the high-speed interface. The low-speed interface that has done the port setup will return a RTS\_Response, either as a discrete message or as part of a CTS. Each CTS sent should be sent from the channel where it is desired that the data be received. An alternative is for all CTS operations to be sent from the same interface, “spoofing” the Source ULA for the operation to make it appear that it was generated by the interface where the data is desired. Using this “spoofing” substitution in combination with a dedicated control channel may also prevent or reduce blocking effects

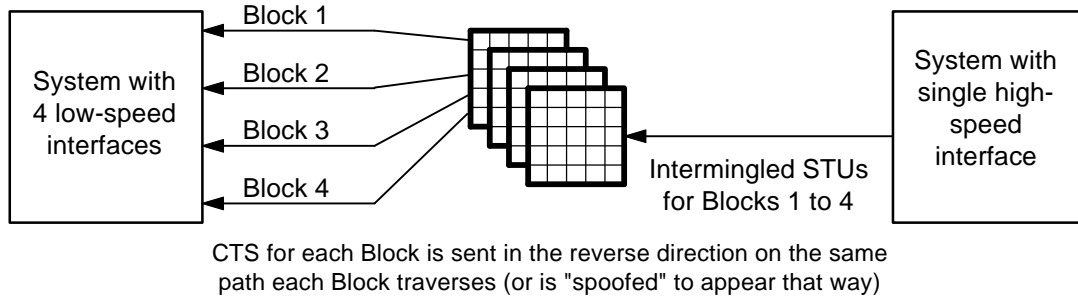
where the underlying physical media suffers from high latency. Subsequent Data Operations may then be done concurrently and will use the Source ULA from in the respective Clear\_To\_Send Operation as the Destination ULA.

Many-to-many striping is the combination of the one-to-many and many-to-one striping. The system receiving data indicates his desire to receive in a striped fashion by issuing multiple CTS operations with differing return Source ULAs. The system sending data chooses to stripe by sending from multiple interfaces that are capable of reaching the proper destination ULA.

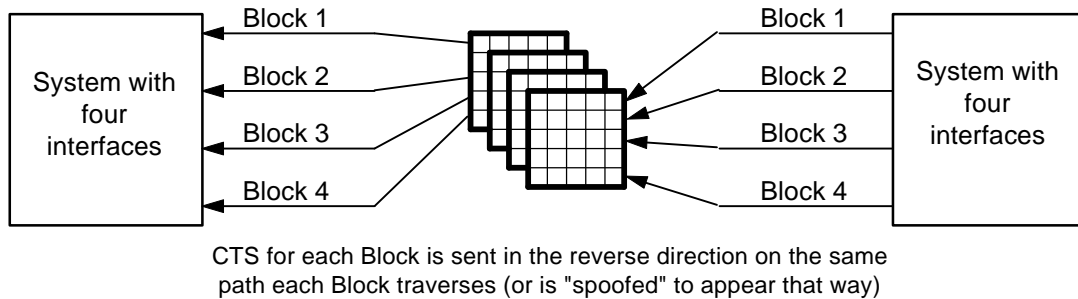
#### B.4 Many-to-many striping



**Figure B.1 – Many-to-one striping**



**Figure B.2 – One-to-many striping**



**Figure B.3 – Many-to-many striping**

## Annex C (informative)

### Scheduled Transfer example

*All of annex C is essentially new, hence no margin bars or highlights.*

The following examples demonstrate various aspects of the scheduled transfer mechanism. The examples are intended for generic application, but may contain some implementation specific ideals. The example types contained in this appendix include:

- Detailed simple transfer example (see C.1)
- Persistent memory example (see C.2)
- Transparently translated, striped Ethernet to HIPPI-6400 example (see C.3)
- Translated, striped HIPPI-800 to HIPPI-6400 example (see C.4)

#### C.1 Detailed simple transfer example

The following example demonstrates the basic Scheduled Transfer Operations. Unlike the later examples, all parameter values are shown in this example. Device X is sending a 256 KB file to device Y in this example.

- Fields that contain a \* are unused and transmitted as zeroes.
- All hex values that take up fewer bits than their field have zeroes in the upper bits.
- Table C.1 shows all fields and corresponding values for the numbered Operations shown below. The Operation description and table location are correlated by the numbers in parenthesis.

##### C.1.1 Virtual Connection set up

Device Y initiates a Virtual Connection set up with device X using the Request\_Port Operation. The Virtual Connection set up is a dual direction connection and either device can begin a Scheduled Transfer after set up, independent of which device initiated the Virtual Connection.

#### Y->X (1)

```
Request_Port (
    S_count: SlotsY,
    R_Port: x'0000',
    S_Port: PortY,
    Bufx: x'10' (BufsizeY = 64 KB),
    Offset: KeyY,
    T_len: x'10' (Max-STUY = 64 KB),
    B_num: EtherType)
```

Device Y provides its Key, Port, available slots, Bufsize, Max-STU, and the EtherType for this Virtual Connection to device X. X binds to the well-known Port **x'0000'** which is used for Port set up. Based on the **EtherType** a corresponding Port value is mapped and sent in the Request\_Port\_Response. Device X responds to the request with an Request\_Port\_Response.

#### X->Y (2)

```
Request_Port_Response (
    S_count: SlotsX,
    R_Port: PortY,
    S_Port: PortX,
    Key: KeyY,
    Bufx: x'0E' (BufsizeX = 16 KB),
    Offset: KeyX,
    T_len: x'0E' (Max-STUX = 16 KB))
```

Device X provides its Key, Port, available Slots, Bufsize, and Max-STU to device Y. Y binds to **KeyY** and **PortY** to associate the response with the above Request\_Port (as opposed to other Request\_Port's that Y may have initiated).

A dual direction Virtual Connection has been established and both sides know the other's Key, Port, available Slots, Bufsize, and Max-STU. The Keys are an authentication value which will stop invalid Operations, (e.g., an inadvertent Port\_Teardown Operation might destroy the Virtual Connection and all ongoing transfers.) The Port values are used to map to upper-layer

entities and may be the same mapping used for Internet style Ports. The Slot value gives the destination flow control power over all Operations requiring a Slot in the handling queue. The Max-STU provides a means for the Final Destination and intermediate translators to declare the maximum size of a data payload, independent of the Final Destination's buffer size. The Buftype (buffer size) field provides buffer tiling information to correctly tile the transfer.

Figure C.1 displays the information provided by each device after the Virtual Connection set up.

Virtual Connection values		
<u>Device X</u>	<u>Device Y</u>	<u>Function</u>
SlotsX	SlotsY	op flow control
PortX	PortY	Port binding
KeyX	KeyY	VC binding
BuftypeX	BuftypeY	FD buffer size
Max-STUX	Max-STUY	STU limit

**Figure C.1 – Virtual Connection information exchanged**

### C.1.2 Scheduled Transfer set up

Device X, the Originating Source for this transfer, initiates a 256 KB transfer using the Virtual Connection established above to device Y, the Final Destination.

#### X->Y (3)

```
Request_To_Send (
  Flags: x'002' (Data channel = 2),
  R_Port: PortY,
  S_Port: PortX,
  Key: KeyY,
  S_id: idX,
  T_len: x'40000' (256 KB))
```

Device Y expects the data to be delivered on Data Channel 2 as specified in the Flags parameter. Device Y checks the Originating Source ID (**idX**) with it's list of acceptable devices. Device Y reads the Transfer length (256 KB), agrees to accept the Transfer, and responds with a Request\_To\_Send\_Response.

#### Y->X (4)

```
Request_To_Send_Response (
  S_count: x'11' (128 KB Blocksize),
  R_Port: PortX,
  S_Port: PortY,
  Key: KeyX,
  R_id: idX,
  S_id: idY,
  T_len: x'40000' (256 KB))
```

Device Y selects a Blocksize of  $2^{17}$  (128 KB). The 256 KB Transfer will consist of two 128 KB Blocks. Device Y assigns an ID (**idY**) which it can use to quickly index to the correct transfer.

All further Scheduled Transfer Operations will continue to use the appropriate binding information: Key, Ports, and the ID's but they are not shown in the remaining steps of this example. See table C.1 for the exhaustive list of parameters for each operation.

### C.1.3 Block 0 transfer

Device Y sends a Clear\_To\_Send Operation once it has finished allocating resources for the Transfer.

#### Y->X (5)

```
Clear_To_Send (
  S_count: x'11' (128 KB Blocksize),
  Bufx: Bufval,
  Offset: x'0',
  B_num: x'0',
  OS_Offset: x'0')
```

The Clear\_To\_Send relays the Destination's starting buffer index (Bufx), the Initial Offset, and the Offset for this Block. The Offset and Initial Offset (carried in the OS\_Offset field) should be the same for the first block of a Transfer. The Initial Offset allows the Originating Source to accurately section its side of the Transfer when the first Block of a Transfer is not sent first.

The Scheduled Transfer Operations concentrate the required complexity on the Originating Source (device X in this case), rather than imposing maximum surprise on the Final Destination.

In this simple example, an efficient tiling between the two nodes is apparent. Device X can send eight 16 KB STU's to correctly fill a Block in

device Y. However, HIPPI-ST accommodates a large range of implementations by allowing any number and size of STU's to fill a block with the following requirements: an STU may not exceed the Max-STU size, 64 KB in this case; and an STU may not be sent to the Final Destination that would overrun a Buffer boundary, Block boundary, or Transfer boundary. For the first Block, the Originating Source, limited by its own buffer size, sends eight 16 KB STU's.

#### X->Y (6-13 are all Data Operations)

#	Flags	S_count	Bufx	Offset	B_num
(6)	x'202'	x'0'	Bufval	x'0'	x'0'
(7)	x'202'	x'1'	Bufval	x'4000'	x'0'
(8)	x'202'	x'2'	Bufval	x'8000'	x'0'
(9)	x'202'	x'3'	Bufval	x'C000'	x'0'
(10)	x'202'	x'4'	Bufval+1	x'0'	x'0'
(11)	x'202'	x'5'	Bufval+1	x'4000'	x'0'
(12)	x'202'	x'6'	Bufval+1	x'8000'	x'0'
(13)	x'00A'	x'7'	Bufval+1	x'C000'	x'0'

As can be seen in this example, the Originating Source does most of the work to align the Buffer regions in the Destination. Figure C.2 shows how the first Block fits into Y's Bufx regions. The last STU of the Block clears the Silent bit so that a Slot resource will be allocated for this STU (informing the ULP of Block reception).

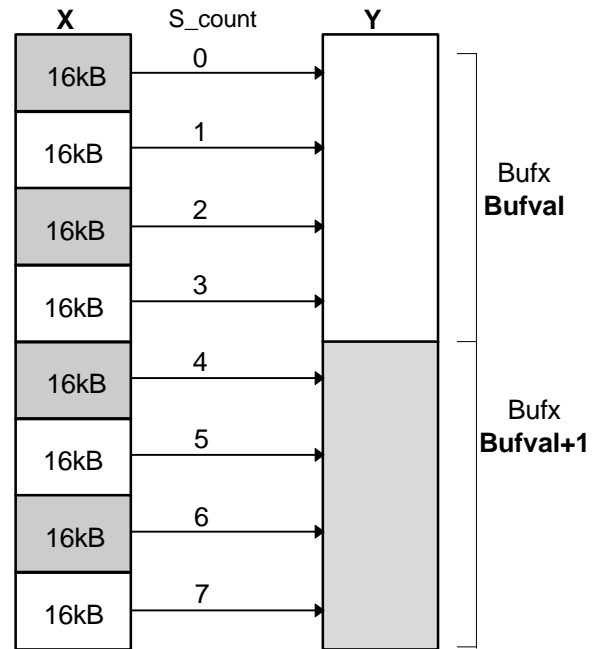


Figure C.2 – Block 0 buffer tiling

#### C.1.4 Block 1 Clear\_To\_Send

The second Clear\_To\_Send allows device X to begin sending the second Block. Because the two Clear\_To\_Sends contain no overlapping buffer regions (in this example), they could have been issued one after another. The second Clear\_To\_Send could also use the same Bufx as the first one, but only if it waits for the first Block to complete before issuing the second Clear\_To\_Send.

*Open Issue - Should the second Block start at a completely different Bufx range, e.g., Bufval+50, to reinforce the idea of each Block being independent of each other?*

#### Y->X (14)

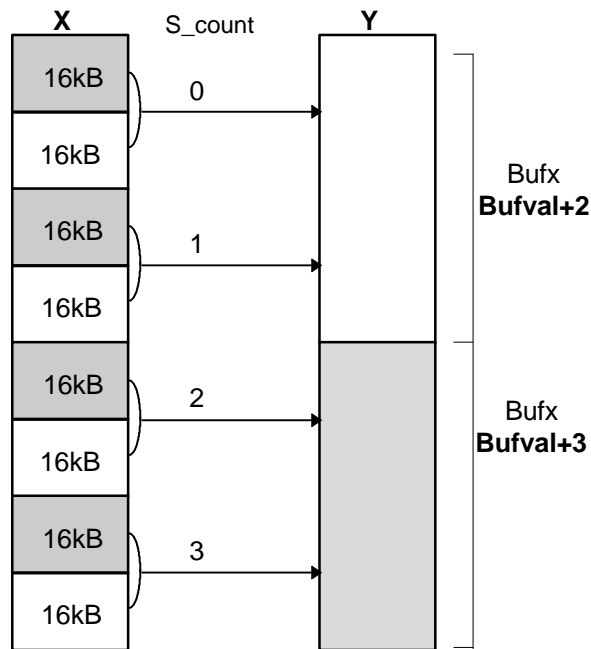
```
Clear_To_Send (
    S_count: x'11' (128 KB Blocksize),
    Bufx: Bufval+2,
    Offset: x'0',
    B_num: x'1',
    OS_Offset: x'0')
```

For the second Block, the Originating Source magically grows a buffer gather mechanism which can pull two buffers at a time. The resulting Data operations appear below.

**X->Y (15-18 are all Data Operations)**

#	Flags	S_count	Bufx	Offset	B_num
(15)	x'202'	x'0'	Bufval+2	x'0'	x'1'
(16)	x'202'	x'1'	Bufval+2	x'8000'	x'1'
(17)	x'202'	x'2'	Bufval+3	x'0'	x'1'
(18)	x'08A'	x'3'	Bufval+3	x'8000'	x'1'

Figure C.3 shows the resulting tiling.



**Figure C.3 – Block 1 buffer tiling**

The gather mechanism decreases the number of STU's sent and nominally improves performance.

*Open Issue - Text about using separate Final Destination ID value for each Block, e.g., second Block uses idY2, goes here.*

The last Data Operation contains the Send\_State flag which triggers device Y to update the Slot value and acknowledge the last Block received. Each of the Data Operations carries a synchronization value in the T\_len parameter which the **Request\_State\_Response** will echo so device X can synchronize it's Slot parameter with the number of outstanding Operations.

**Y->X (19)**

**Request\_State\_Response** (  
 S\_count: **C-SlotsY**,  
 Offset: **x'1'**,  
 T\_len: **SyncX**,  
 B\_num: **x'1'**)

The Request\_State\_Response contains an updated Slot value, called **C-SlotsY**, in S\_count. The synchronization value is contained in T\_len; the Block number acknowledges this Block and all lower numbered Blocks. The Offset contains the highest received Block number in a contiguous set from the first Block. The value in B\_num acknowledges the Block that had the Send\_State flag set.

**C.1.5 Ending the Virtual Connection**

Either side can terminate the Virtual Connection and free all of the resources associated with the Virtual Connection.

**Y->X (20)**

Port\_Teardown (  
 R\_Port: **PortX**,  
 S\_Port: **PortY**,  
 Key: **KeyX**)

**X->Y (21)**

Port\_Teardown\_ACK (  
 R\_Port: **PortY**,  
 S\_Port: **PortX**,  
 Key: **KeyY**)

**Y->X (22)**

Port\_Teardown\_Complete (  
 R\_Port: **PortX**,  
 S\_Port: **PortY**,  
 Key: **KeyX**)

The Port\_Teardown is a three-way handshake that decreases timeout dependency for releasing resources. The device sending the Port\_Teardown\_ACK can release all of the resources upon reception of the Port\_Teardown\_Complete.

Table C.1 – Scheduled Transfer example summary

Operation	Op	Flags	S_count	R_Port	S_Port	Key	R_id	S_id	Bufx	Offset	T_len	B_num
<b>(1) Request_Port (Y-&gt;X)</b>	<b>x'01'</b>		<b>SlotsY</b>	<b>x'0000'</b>	<b>PortY</b>	*	*	*	<b>x'10'</b> <b>(64 KB)</b>	<b>KeyY</b>	<b>x'10'</b> <b>(64 KB)</b>	<b>Ether-Type</b>
(2) RQP_Response (X->Y)	x'02'		SlotsX	PortY	PortX	KeyY	*	*	x'0E' (16 KB)	KeyX	x'0E' (16 KB)	*
(3) Request_To_Send (X->Y)	x'06'	D2	*	PortY	PortX	KeyY	*	idX	*	*	x'40000'	
<b>(4) RTS_Response (Y-&gt;X)</b>	<b>x'07'</b>		<b>x'11'</b>	<b>PortX</b>	<b>PortY</b>	<b>KeyX</b>	<b>idX</b>	<b>idY</b>	*	*	*	*
<b>(5) Clear_To_Send (Y-&gt;X)</b>	<b>x'0A'</b>		<b>x'11'</b>	<b>PortX</b>	<b>PortY</b>	<b>KeyX</b>	<b>idX</b>	<b>idY</b>	<b>Bufval</b>	<b>x'0'</b>	*	<b>x'0'</b>
(6) Data (X->Y)	x'0B'	T, D2	0	PortY	PortX	KeyY	idY	idX	Bufval	x'0'	Syncval	x'0'
(7) Data (X->Y)	x'0B'	T, D2	1	PortY	PortX	KeyY	idY	idX	Bufval	x'4000'	Syncval	x'0'
(8) Data (X->Y)	x'0B'	T, D2	2	PortY	PortX	KeyY	idY	idX	Bufval	x'8000'	Syncval	x'0'
(9) Data (X->Y)	x'0B'	T, D2	3	PortY	PortX	KeyY	idY	idX	Bufval	x'C000'	Syncval	x'0'
(10) Data (X->Y)	x'0B'	T, D2	4	PortY	PortX	KeyY	idY	idX	Bufval+1	x'0000'	Syncval	x'0'
(11) Data (X->Y)	x'0B'	T, D2	5	PortY	PortX	KeyY	idY	idX	Bufval+1	x'4000'	Syncval	x'0'
(12) Data (X->Y)	x'0B'	T, D2	6	PortY	PortX	KeyY	idY	idX	Bufval+1	x'8000'	Syncval	x'0'
(13) Data (X->Y)	x'0B'	L, D2	7	PortY	PortX	KeyY	idY	idX	Bufval+1	x'C000'	Syncval	x'0'
<b>(14) Clear_To_Send (Y-&gt;X)</b>	<b>x'0A'</b>		<b>x'11'</b>	<b>PortX</b>	<b>PortY</b>	<b>KeyX</b>	<b>idX</b>	<b>idY2</b>	<b>Bufval+2</b>	<b>x'0'</b>		<b>x'1'</b>
(15) Data (X->Y)	x'0B'	T, D2	0	PortY	PortX	KeyY	idY2	idX	Bufval+2	x'0'	Syncval	x'1'
(16) Data (X->Y)	x'0B'	T, D2	1	PortY	PortX	KeyY	idY2	idX	Bufval+2	x'8000'	Syncval	x'1'
(17) Data (X->Y)	x'0B'	T, D2	2	PortY	PortX	KeyY	idY2	idX	Bufval+3	x'0'	Syncval	x'1'
(18) Data (X->Y)	x'0B'	S,L, D2	3	PortY	PortX	KeyY	idY2	idX	Bufval+3	x'8000'	SyncX	x'1'
<b>(19) Request_State_Response (Y-&gt;X)</b>	<b>x'0E'</b>		<b>C-SlotsY</b>	<b>PortX</b>	<b>PortY</b>	<b>KeyX</b>	<b>idX</b>	<b>idY2</b>	*	<b>x'1'</b>	<b>SyncX</b>	<b>x'1'</b>
<b>(20) Port_Teardown (Y-&gt;X)</b>	<b>x'03'</b>	*	*	<b>PortX</b>	<b>PortY</b>	<b>KeyX</b>	*	*	*	*	*	*
(21) Port_Teardown_ACK (X->Y)	x'04'	*	*	PortY	PortX	KeyY	*	*	*	*	*	*
<b>(22) Port_Teardown_Complete (Y-&gt;X)</b>	<b>x'05'</b>	*	*	<b>PortX</b>	<b>PortY</b>	<b>KeyX</b>	*	*	*	*	*	*
NOTE – Operations from X to Y are shown in plain text; Operations from Y to X are shown in bold italic.												

## C.2 Persistent memory example

Though not detailing a shared memory structure, the following example presents the set up of a buffer region that handles both puts and gets, a building block for shared memory operations. Figure C.4 shows the devices and Bufx range used in this example.

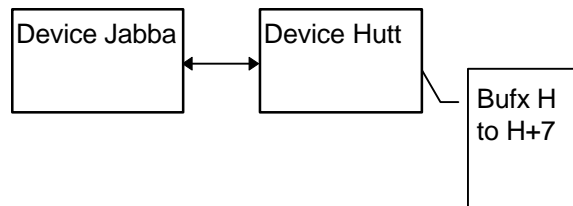


Figure C.4 – Persistent memory setup

### C.2.1 Set up

The following commands set up the Virtual Connection:

#### Jabba → Hutt

```
Request_Port (
  Slots: SlotsJabba,
  R_Port: SMPort,
  S_Port: SMPort,
  Bufsize: x'0E' (16 KB),
  XKey: KeyJabba,
  Max-STU: x'0E' (16 KB),
  Ethertype: Local)
```

#### Hutt → Jabba

```
Request_Port_Response (
  Flags: x'010' (Persistent),
  Slots: SlotsHutt,
  R_Port: SMPort,
  S_Port: SMPort,
  Key: KeyJabba,
  Bufsize: x'0D' (8 KB),
  XKey: KeyHutt,
  Max-STU: x'0D' (8 KB))
```

Device Hutt flags that he accepts Persistent Request\_To\_Send transfers. The rest of the parameters are summarized in table C.2

Table C.2 – Virtual Connection parameters

Parameter	Jabba	Hutt
Slots	SlotsJabba	SlotsHutt
Ports	SMPort	SMPort
Keys	KeyJabba	KeyHutt
Bufsize	16 KB	8 KB
Max-STU	16 KB	8 KB

The two devices use a common 64 KB network memory line size. Device Jabba would like to read and write memory on device Hutt of this size and issues an Request\_To\_Send.

#### Jabba → Hutt

```
Request_To_Send (
  Flags: x'011' (Persistent, Channel 1),
  S_id: idJ,
  T_len: x'10000' (64 KB))
```

Parameters in the optional 32 bytes of the Request\_To\_Send Operation may identify the exact memory location device Jabba is asking to use. Device Hutt maps the memory to a Bufx and returns the Bufx to Jabba in a Clear\_To\_Send Operation.

#### Hutt → Jabba

```
Clear_To_Send (
  Blocksize: x'10' (64 KB),
  R_id: idJ,
  S_id: idtH,
  Bufx: BufxH,
  Offset: x'0',
  T_len: x'10000' (64 MB),
  B_num: x'0',
  I_Offset: x'0')
```

Device Hutt has pinned down the persistent memory page allowing device Jabba to either read (with an Request\_To\_Receive Operation) or write (with a Data Operation). Note that with Persistent = 1, Data Operations can be issued without being having to be preceded by a Clear\_To\_Send Operation from the other end.



### C.2.2 Reading

Device Jabba would first like to read the established persistent memory Block prior to changing values and initiates a Request\_To\_Send to device Hutt.

#### Jabba → Hutt

```
Request_To_Send (
  Flags: x'011' (Persistent, Channel 1),
  R_id: idH,
  S_id: idJ,
  Bufx: BufxH,
  Offset: x'0',
  T_len: x'10000' (64 KB),
  OS_Bufx: BufxJ,
  OS_Offset: x'0')
```

Hutt receives the Request\_To\_Receive Operation and responds with a number of Data Operations. Each STU is 8 KB, limited by Hutt's send buffer mechanism. Hence, in this example eight Data Operations will be used, i.e., 64 KB Block read divided by the STU size. Only the first Data Operation is shown in detail, parameter differences in the following ones are listed below it.

#### Hutt → Jabba

```
#1, Data (
  Flags: x'201',
  S_count: x'0',
  R_id: idJ,
  S_id: idH,
  Bufx: BufxJ,
  Offset: x'0',
  B_num: x'0')
```

#	Flags	S_count	Bufx	Offset
2	x'201'	x'1'	BufxJ	x'2000'
3	x'201'	x'2'	BufxJ+1	x'0'
4	x'201'	x'3'	BufxJ+1	x'2000'
5	x'201'	x'4'	BufxJ+2	x'0'
6	x'201'	x'5'	BufxJ+2	x'2000'
7	x'201'	x'6'	BufxJ+3	x'0'
8	x'209'	x'7'	BufxJ+3	x'2000'

Device Jabba has successfully read the contents of the memory location set up on device Hutt. Hutt could have set the Send\_State bit in the last Data Operation if it cared to receive an

acknowledgment that device Jabba completed the read.

### C.2.3 Writing

Next, Jabba would like to change the memory location on Hutt and sends a Data Operation without the baggage of Block flow control or other setup (after the initial Request\_To\_Send, Clear\_To\_Send). Each STU is 8 KB to fill one of device Hutt's buffers. Only the first Data Operation is shown in detail, parameter differences in the ones following are listed below it.

#### Jabba → Hutt

```
#1, Data (
  Flags: x'201',
  S_count: x'0',
  R_id: idH,
  S_id: idJ,
  Bufx: BufxH,
  Offset: x'0',
  B_num: x'0')
```

#	Flags	S_count	Bufx	Offset
2	x'201'	x'1'	BufxH	x'2000'
3	x'201'	x'2'	BufxH+1	x'0'
4	x'201'	x'3'	BufxH+1	x'2000'
5	x'201'	x'4'	BufxH+2	x'0'
6	x'201'	x'5'	BufxH+2	x'2000'
7	x'201'	x'6'	BufxH+3	x'0'
8	x'089'	x'7'	BufxH+3	x'2000'

The last Data Operation sets the Send\_State bit to confirm reception of the write. Silent = 0 is used to alert device Hutt that the location has been written.

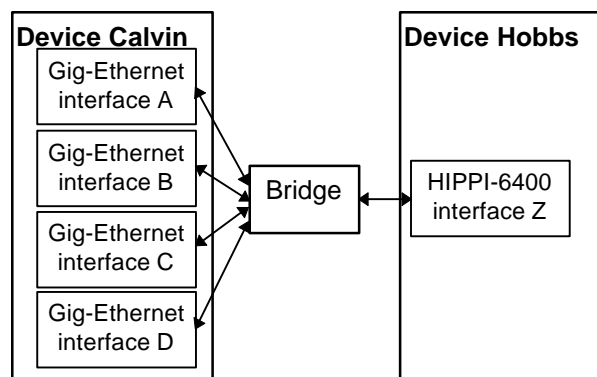
Device Jabba could perform multiple read and writes into this location with very little overhead. Many of these locations could also be set up on device Hutt along with multiple outside devices accessing the same locations.

### C.2.4 Closing the persistent memory

Either side can close the persistent memory location by issuing an END operation. Port\_Teardown Operations will also close the memory location.

### C.3 Translated, striped Ethernet to HIPPI-6400 example

A simple, though common application of striping takes multiple Gigabit Ethernet links and aggregates them onto HIPPI-6400. The duties of both end devices and the translator are detailed in this example. Figure C.5 shows the network topology for this example.



**Figure C.5 – Network topology**

Calvin and Hobbs are two end devices that need to communicate large amounts of data. Unfortunately, Calvin only speaks Gigabit Ethernet and Hobbs only knows HIPPI-6400. A transparent, Gigabit-Ethernet-switching translator is placed between them to break the MAC protocol barrier. To enable high throughput direct memory access (DMA) transfers, the devices use the HIPPI-ST protocol.

The following example demonstrates data transfers moving in both directions between Calvin and Hobbs.

#### C.3.1 Virtual Connection Setup

Hobbs initiates the Virtual Connection setup to Calvin's interface **A** ULA address (found using ARP for Calvin). Interface **A** will act as the control channel interface as well as one of the Data pipes. A non-Data control interface could also be used to reduce any contention between Control and Data operations on interface **A**. The translator will take the HIPPI-6400 Request\_Port Operation, look up the correct Gigabit Ethernet port for **A**, reformat the HIPPI-6400 message as a Gigabit Ethernet packet, and forward the packet to Gigabit Ethernet interface **A**.

Hobbs inserts a zero into the Max-Blocksize (carried in the OS\_Bufx) parameter. The translator checks its available buffers for this port and inserts a Max-Blocksize value of **x'13'** (512 KB) into the Request\_Port before forwarding it to Calvin.

**Z → A**

**Request\_Port (**

Flags: **x'400'** (*Out\_of\_Order*),  
Slots: **SlotsHobbs**,  
R\_Port: **x'0000'**,  
S\_Port: **PortHobbs**,  
Bufsize: **x'0E'** (16 KB),  
XKey: **KeyHobbs**,  
Max-STU: **x'0E'** (16 KB),  
Ethertype: **Ethertype**,  
Max-Blocksize: **x'0' → x'13'** (512 KB))

*Open Issue – Aghh! What's this new parameter? The Max-Blocksize states the maximum data that can be stored and forwarded through the translator without a strict flow control mechanism running between the translator and both ends. This means that the translator can have less complicated ST protocol hardware, but requires large, high-speed buffers.*

Calvin processes the request and sends a Request\_Port\_Response from interface **A**. The translator receives the Request\_Port\_Response, converts the Gigabit Ethernet packet format to a HIPPI-6400 message, and transmits the message on Virtual Channel 0 (because it's a Control Operation). Again, the translator overwrites the Max-Blocksize value.

**A → Z**

**Request\_Port\_Response (**

Flags: **O**,  
Slots: **SlotsCalvin**,  
R\_Port: **PortHobbs**,  
S\_Port: **PortCalvin**,  
Key: **KeyHobbs**,  
Bufsize: **x'0C'** (4 KB),  
XKey: **KeyCalvin**,  
Max-STU: **x'0A'** (1 KB),  
Max-Blocksize: **x'0' → x'13'** (512 KB))

Calvin has a 4 KB Bufsize and selects a 1 KB Max-STU size to accommodate the 1500 byte size constraint of Gigabit Ethernet packets. The Virtual Connection is established and both devices have indicated the ability to receive Blocks out of order (required for translated

striping). Table C.3 summarizes the parameters exchanged.

**Table C.3 – Virtual Connection parameters**

Parameter	Calvin	Hobbs
<b>Slots</b>	SlotsCalvin	SlotsHobbs
<b>Ports</b>	PortCalvin	PortHobbs
<b>Keys</b>	KeyCalvin	KeyHobbs
<b>Bufsize</b>	4 KB	16 KB
<b>Max-STU</b>	1 KB	16 KB
<b>Max-Blocksize</b>	512 KB	512 KB

### C.3.2 Sending to Hobbs

Calvin needs to send a one megabyte Transfer to Hobbs.

**A → Z**

**Request\_To\_Send (**  
 Flags: **x'001'**,  
 Outstanding\_Blocks: **x'10'**,  
 S\_id: **idt1C**,  
 T\_len: **x'100000'** (1 MB))

*Open Issue – Note the new parameter in S\_count called “Outstanding Blocks”. The Orig. Source requires a means to “advise” the Final Destination as to Blocksize selection to keep the Orig. Source’s interfaces busy.*

Calvin selects to send Data Operations through Virtual Channel 1. Calvin sets his preferred number of outstanding blocks to 16, hopefully, keeping four blocks outstanding for each of the four interfaces. The Outstanding Blocks parameter is only a suggestion and may be ignored by Hobbs.

A total of 16 Clear\_To\_Send Operations are sent from **Z** to **A**. Only the first Clear\_To\_Send is shown in detail, parameter differences in the ones following are listed below it.

**Z → A**

**#1, Clear\_To\_Send (**  
 Blocksize: **x'0F'** (32 KB),  
 R\_id: **idt1C**,  
 S\_id: **idt1H**,  
 Bufx: **Hobbs0**,  
 Offset: **x'0'**,  
 B\_num: **x'0'**)

#	Bufx	B_num
2	Hobbs2	x'1'
3	Hobbs4	x'2'
4	Hobbs6	x'3'
5	Hobbs8	x'4'
6	Hobbs10	x'5'
7	Hobbs12	x'6'
8	Hobbs14	x'7'
9	Hobbs16	x'8'
10	Hobbs18	x'9'
11	Hobbs20	x'A'
12	Hobbs22	x'B'
13	Hobbs24	x'C'
14	Hobbs26	x'D'
15	Hobbs28	x'E'
16	Hobbs30	x'F'

Hobbs selects a Blocksize of 32 KB based on his buffer size, the Max-Blocksize, and his ability to handle interrupts. Hobbs sends 16 Clear\_To\_Sends to interface **A** meeting the requested number of outstanding Blocks, while staying under the Max-Blocksize.

Calvin begins receiving the Clear\_To\_Send requests through interface **A** and in a round-robin fashion farms Data Operation requests to each of the four Gigabit Ethernet interfaces. Only the first Data Operation is shown in detail, parameter differences in the ones following are listed below it.

**A → Z**

**Data (**  
 Flags: **x'201'**,  
 S\_count: **x'0'**,  
 R\_id: **idt1H**,  
 S\_id: **idt1C**,  
 Bufx: **Hobbs0**,  
 Offset: **x'0'**,  
 B\_num: **x'0'**)

#### **B → Z**

```
Data (
  Flags: x'201',
  S_count: x'0',
  R_id: idt1H,
  S_id: idt1C,
  Bufx: Hobbs2,
  Offset: x'0',
  B_num: x'1')
```

#### **C → Z**

```
Data (
  Flags: x'201',
  S_count: x'0',
  R_id: idt1H,
  S_id: idt1C,
  Bufx: Hobbs4,
  Offset: x'0',
  B_num: x'2')
```

#### **D → Z**

```
Data (
  Flags: x'201',
  S_count: x'0',
  R_id: idt1H,
  S_id: idt1C,
  Bufx: Hobbs6,
  Offset: x'0',
  B_num: x'3')
```

#### **A → Z**

```
Data (
  Flags: x'201',
  S_count: x'1',
  R_id: idt1H,
  S_id: idt1C,
  Bufx: Hobbs0,
  Offset: x'400' (1 KB),
  B_num: x'0')
```

... (506 1 KB STU's later) ...

#### **D → Z**

```
Data (
  Flags: x'089',
  S_count: x'31',
  R_id: idt1H,
  S_id: idt1C,
  Bufx: Hobbs31,
  Offset: x'3C00' (15 KB),
  B_num: x'15')
```

Each of the Data Operations will contain 1 KB of STU for the simplest tiling scheme (due to the

Gigabit Ethernet packet size restriction). The flag bits carry the data channel value so the translator can correctly place each STU on the correct Virtual Channel. All STU's are marked silent except for the last of each Block which requests a Request\_State\_Response acknowledgment. The STU count will range from 0 to 31 for each block (32 each 1 KB STU's). The last Data operation listed above shows interface **D** sending the last STU of Block 15 (halfway through the transfer). The order that the Blocks are sent may not match the order shown above which is why Hobbs must accept out of order Blocks. The translator is required to switch between the incoming Gigabit Ethernet interfaces onto VC1 of the HIPPI-6400 interface. The aggregate bandwidth of the striped Gigabit Ethernet sits around half the bandwidth of the HIPPI-6400 link and things should flow smoothly, (note the translator's capability to buffer all the outstanding blocks in case a lack of credits appears on the HIPPI-6400 link).

As groups of 32 1 KB STU's fill in each Block in the HIPPI-6400 receive buffers, new blocks may be cleared for transmission (as long as the 512 KB translator buffer is not overrun). Blocks in error may be resent according to normal HIPPI-6400 retransmission methods.

Hobbs will clear another 16 Blocks for transmission to complete the megabyte Transfer. Acknowledgments are sent in the normal manner.

### **C.3.3 Sending to Calvin**

Hobbs received the first transfer and after changing a bit, sends the megabyte of data back to Calvin.

#### **Z → A**

```
Request_To_Send (
  Flags: x'002',
  Outstanding_Blocks: x'0',
  S_id: idt2H,
  T_len: x'100000' (1 MB))
```

Hobbs plans to send on Data Channel 2 and sets the outstanding Blocks to zero showing that Hobbs doesn't care. This will allow Calvin to select a Blocksize that is not constrained by a set number of Blocks, (but still by the Max-Blocksize).

**A → Z**

```
Clear_To_Send (
  Blocksize: x'11' (128 KB),
  R_id: idt2H,
  S_id: idt2C,
  Bufx: Calvin0,
  Offset: x'0',
  B_num: x'0')
```

**B → Z**

```
Clear_To_Send (
  Blocksize: x'11' (128 KB),
  R_id: idt2H,
  S_id: idt2C,
  Bufx: Calvin32,
  Offset: x'0',
  B_num: x'1')
```

**C → Z**

```
Clear_To_Send (
  Blocksize: x'11' (128 KB),
  R_id: idt2H,
  S_id: idt2C,
  Bufx: Calvin64,
  Offset: x'0',
  B_num: x'2')
```

**D → Z**

```
Clear_To_Send (
  Blocksize: x'11' (128 KB),
  R_id: idt2H,
  S_id: idt2C,
  Bufx: Calvin96,
  Offset: x'0',
  B_num: x'3')
```

**A → Z**

```
Clear_To_Send (
  Blocksize: x'11' (128 KB),
  R_id: idt2H,
  S_id: idt2C,
  Bufx: Calvin128,
  Offset: x'0',
  B_num: x'4')
```

**B → Z**

```
Clear_To_Send (
  Blocksize: x'11' (128 KB),
  R_id: idt2H,
  S_id: idt2C,
  Bufx: Calvin160,
  Offset: x'0',
  B_num: x'5')
```

**C → Z**

```
Clear_To_Send (
  Blocksize: x'11' (128 KB),
  R_id: idt2H,
  S_id: idt2C,
  Bufx: Calvin192,
  Offset: x'0',
  B_num: x'6')
```

**D → Z**

```
Clear_To_Send (
  Blocksize: x'11' (128 KB),
  R_id: idt2H,
  S_id: idt2C,
  Bufx: Calvin224,
  Offset: x'0',
  B_num: x'7')
```

Each of Calvin's interfaces responds with two **Clear\_To\_Send's** containing its ULA in the Source Address (for the correct return of Data Operations to each interface). The Blocksize of 128 KB means that only four Blocks may be outstanding without overrunning the translator's Max-Blocksize value. Because eight Blocks have been cleared, the translator will queue the last four **Clear\_To\_Send's** until the previous ones have been sent by the translator. Since only a single Block will remain outstanding to any of the Gigabit Ethernet interfaces, some lag time will result in clearing the next Block.

**Z → A**

```
Data(Flags: x'202', S_count: x'0', R_id: idt2C,
S_id: idt2H, Bufx: Hobbs0, Offset: x'0',
B_num: x'0')
```

**Z → B**

```
Data(Flags: x'202', S_count: x'0', R_id: idt2C,
S_id: idt2H, Bufx: Hobbs2, Offset: x'0',
B_num: x'1')
```

**Z → C**

```
Data(Flags: x'202', S_count: x'0', R_id: idt2C,
S_id: idt2H, Bufx: Hobbs4, Offset: x'0',
B_num: x'2')
```

**Z → D**

```
Data(Flags: x'202', S_count: x'0', R_id: idt2C,
S_id: idt2H, Bufx: Hobbs6, Offset: x'0',
B_num: x'3')
```

**Z → A**

**Data**(Flags: **x'202'**, S\_count: **x'1'**, R\_id: **idt2C**, S\_id: **idt2H**, Bufx: **Hobbs0**, Offset: **x'400'** (1 KB), B\_num: **x'0'**)

...(506 1 KB STU's later) ...

**Z → D**

**Data**(Flags: **x'08A'**, S\_count: **x'127'**, R\_id: **idt2C**, S\_id: **idt2H**, Bufx: **Hobbs31**, Offset: **x'1FC00'** (127 KB), B\_num: **x'3'**)

Obeying the Max-STU, Hobbs sends 1024 1 KB STU's, but the larger blocks may cause lulls during the transfer when Hobbs has not received a new **Clear\_To\_Send** to keep things going.

NOTES –

1 – In order for people to actually implement and make things useful, striping needs to be kept simple or it will be the option that doesn't happen. I believe this partially includes somewhat restricting the “applicable in all situations” attributes that are nice, but not always necessary.

2 – Setting Clear\_To\_Send's small enough to keep Blocks in motion without overrunning the Max-Blocksize parameter.

3 – Efficient use of the translator's buffers must be maintained. The Max-Blocksize communicated during port setup must be shared by all transfers. The value is called Max-Blocksize and not intermediate storage size as the host may produce more Clear\_To\_Send's worth in data than the translator can buffer (as long as the Blocksize is not larger than the translator's declared maximum). The translator will queue the Clear\_To\_Send requests forwarding only the number that can be safely handled by the translator. Once the block for an outstanding Clear\_To\_Send is sent, the next Clear\_To\_Send in the queue may proceed to the end host.

4 – This is all new and the only way the author could think of to enable a now semi-transparent translator between ST devices.

5 – A less complex method can be easily applied to simple N to N stripes across the same media. Neither of the suggested new parameters are required for direct N to N striping.